

✓ 315.930  
Studia

13 a

20

20  
1985

# Scientiarum Mathematicarum Hungarica

EDITOR-IN-CHIEF

A. HAJNAL

DEPUTY EDITOR-IN-CHIEF

E. T. SCHMIDT

EDITORIAL BOARD

P. BOD, E. CSÁKI, Á. CSÁSZÁR, I. CSISZÁR, Á. ELBERT

L. FEJES TÓTH, G. HALÁSZ, I. JUHÁSZ, G. KATONA

O. STEINFELD, J. SZABADOS, D. SZÁSZ

E. SZEMERÉDI, G. TUSNÁDY, I. VINCZE, R. WIEGANDT

VOLUME 20  
NUMBERS 1—4  
1985



AKADÉMIAI KIADÓ, BUDAPEST

# STUDIA SCIENTIARUM MATHEMATICARUM HUNGARICA

A QUARTERLY OF THE HUNGARIAN  
ACADEMY OF SCIENCES

---

*Studia Scientiarum Mathematicarum Hungarica* publishes original papers on mathematics mainly in English, but also in German, French and Russian.

*Studia Scientiarum Mathematicarum Hungarica* is published in yearly volumes of four issues (mostly double numbers published semiannually) by

AKADÉMIAI KIADÓ

Publishing House of the Hungarian Academy of Sciences  
H-1054 Budapest, Alkotmány u. 21.

Manuscripts and editorial correspondence should be addressed to

J. Merza  
Managing Editor  
P.O. Box 127  
H-1364 Budapest

## *Subscription information*

Orders should be addressed to

KULTURA Foreign Trading Company  
P.O. Box 149  
H-1389 Budapest

or to its representatives abroad.

# Studia Scientiarum Mathematicarum Hungarica

Editor-in-Chief

A. Hajnal

Deputy Editor-in-Chief

E. T. Schmidt

Editorial Board

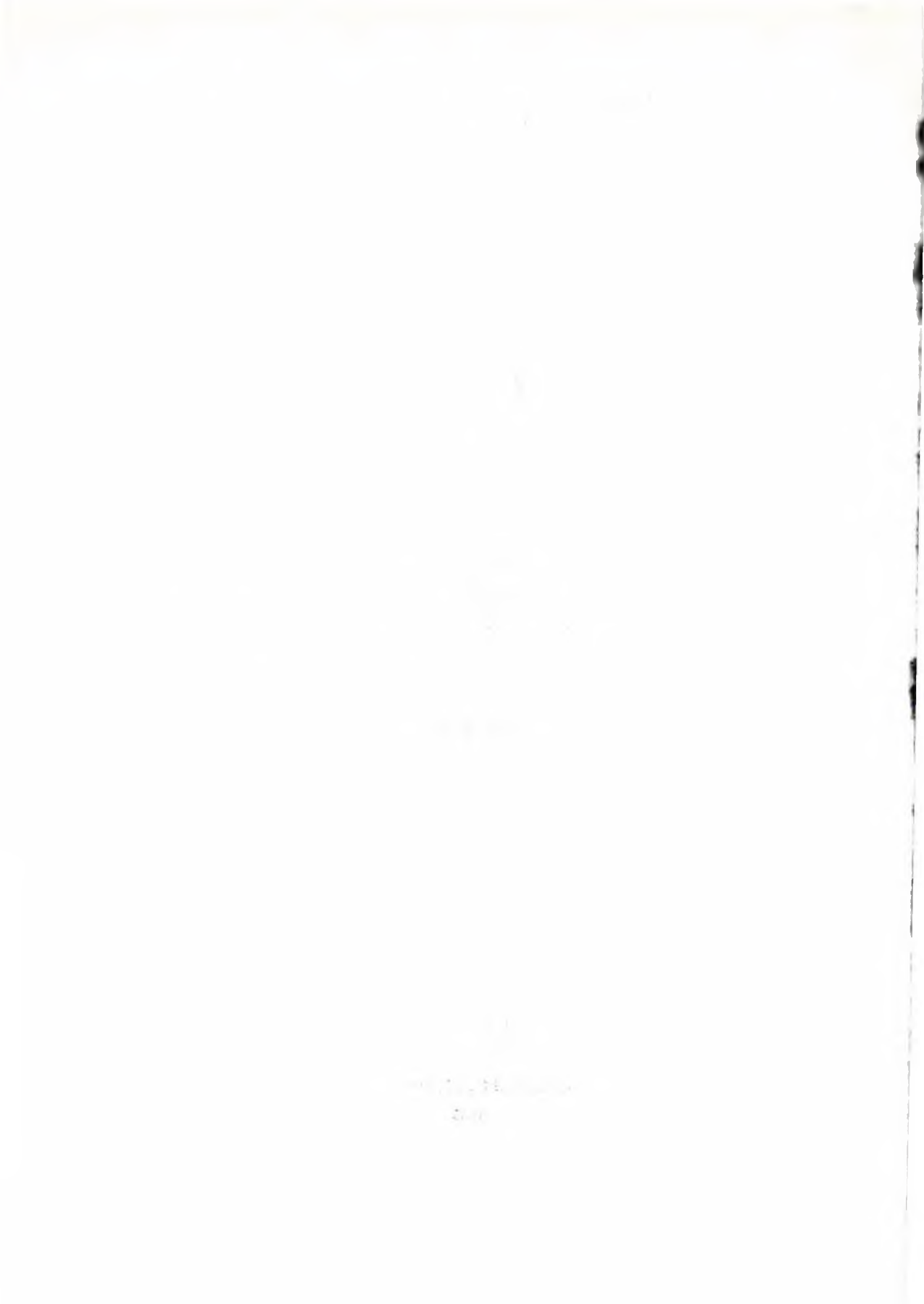
P. Bod, E. Csáki, Á. Császár, I. Csiszár, Á. Elbert, L. Fejes Tóth  
G. Halász, I. Juhász, G. Katona, O. Steinfeld, J. Szabados, D. Szász  
E. Szemerédi, G. Tusnády, I. Vincze, R. Wiegandt

Volume 20



Akadémiai Kiadó, Budapest

1985





## CONTENTS

<i>Ahmad, M.</i> , Estimation of the parameters of Burr distribution based on order statistics . . . . .	133
<i>Bayasgalan, Ts.</i> , О фундаментальной приводимости самосопряженных и унитарных операторов в пространствах с индефинитной метрикой . . . . .	313
<i>Beazer, R.</i> , Congruence uniform algebras with pseudocomplementation . . . . .	43
<i>Beck, J.</i> , Remarks on combinatorial geometry I . . . . .	249
<i>Bell, H. E.</i> , On commutativity of quasi-commutative rings . . . . .	21
<i>Berg, G.</i> , Steinness and the vanishing of cohomology . . . . .	193
<i>Bihari, I.</i> , An asymptotic statement concerning the solutions of the differential equation $x'' + a(t)x = 0$ . . . . .	11
<i>Bihari, I.</i> , Note to an extension of a Sturmian comparison theorem . . . . .	15
<i>Blasco, J. L.</i> , Complete bases in topological spaces . . . . .	49
<i>Book review</i> . . . . .	481
<i>Chiang, C.-Y. and Puri, M. L.</i> , Tests of subhypotheses in linear regression based on rank-order estimates . . . . .	237
<i>Colbourn, C. J. and Rosa, A.</i> , Indecomposable triple systems with $\lambda=4$ . . . . .	139
<i>Damaschke, P. and Stern, M.</i> , A characterization of generalized matroid lattices . . . . .	425
<i>Darbinyan, S. Kh.</i> , Пандикличность орграфов при условии Мейнила . . . . .	95
<i>Duggal, B. P.</i> , On the spectrum of a class of integral transforms II . . . . .	451
<i>Duggal, B. P.</i> , On the spectrum of a class of integral transforms III. An application . . . . .	461
<i>Eigen, S. J.</i> , Putting convergent sequences into measurable sets . . . . .	411
<i>Fejes Tóth, L.</i> , Packing of homothetic discs of $n$ different sizes . . . . .	217
<i>Fialowski, A.</i> , On the deformations of $L_1$ . . . . .	433
<i>Florian, A.</i> , On compact packing of circles . . . . .	473
<i>Florian, A. and Groemer, H.</i> , Two remarks on the permeability of layers of convex bodies . . . . .	259
<i>Frankl, P.</i> , Bounding the size of a family knowing the cardinality of differences . . . . .	33
<i>Gaál, I.</i> , Norm form equations with several dominating variables and explicit lower bounds for inhomogeneous linear forms with algebraic coefficients II. . . . .	333
<i>Grimmett, G. R.</i> , The largest components in a random lattice . . . . .	325
<i>Groemer, H. and Florian, A.</i> , Two remarks on the permeability of layers of convex bodies . . . . .	259
<i>Grossman, E. H.</i> , Number bases in quadratic fields . . . . .	55
<i>Gut, A.</i> , On the law of the iterated logarithm for randomly indexed partial sums with two applications . . . . .	63
<i>Harman, G., Pintz, J. and Wolke, D.</i> , A note on the Möbius and Liouville functions . . . . .	295
<i>Huijsmans, C. B.</i> , An inequality in complex Riesz algebras . . . . .	29
<i>Huynh, Dinh Van</i> , On rings with modified chain conditions . . . . .	59
<i>Imhof, J. P.</i> , On Brownian bridge and excursion . . . . .	1
<i>Isac, G.</i> , Branches continues de vecteurs propres généralisés. Applications aux équations de coïncidences . . . . .	155
<i>Ivič, A.</i> , The distribution of primitive abundant numbers . . . . .	183
<i>Ivič, A.</i> , On squarefree numbers with restricted prime factors . . . . .	189
<i>Khan, L. A.</i> , Separability in the uniform topology . . . . .	407
<i>Khare, S. S.</i> , Reduction of equivariant bordism groups . . . . .	213
<i>Khoi, Trinh Dang</i> , Строго наследственные радикалы в классе всех топологических колец . . . . .	37
<i>Kiss, P.</i> , Differences of the terms of linear recurrences . . . . .	285
<i>Komornik, V.</i> , On the eigenfunctions of first- and second-order differential operators . . . . .	275
<i>Lénárd, M.</i> , Spline interpolation in two variables . . . . .	145
<i>Löfström, J.</i> , Best approximation in $L_p(w)$ by algebraic polynomials . . . . .	375

<i>Marcus, F.</i> , Sur les surfaces à groupes continus $G_2$ de similitude projectives en elles-mêmes et sur les surfaces complexes .....	267
<i>Miller, H. I.</i> and <i>Xenikakis, P. J.</i> , Some results on the Cantor set .....	309
<i>Moór, A.</i> und <i>Nadj, D. F.</i> , Über die autoparallele Abweichung von Finsler—Otsukischen Räumen und Anwendungen in Räumen mit speziellen $P$ -Tensoren .....	395
<i>Nadj, D. F.</i> und <i>Moór, A.</i> , Über die autoparallele Abweichung von Finsler—Otsukischen Räumen und Anwendungen in Räumen mit speziellen $P$ -Tensoren .....	395
<i>Panny, W.</i> and <i>Prodinger, H.</i> , The expected height of paths for several notions of height .....	119
<i>Perelli, A.</i> and <i>Salerno, S.</i> , On $2k$ -dimensional density estimates .....	345
<i>Pintz, J.</i> , <i>Harman, G.</i> and <i>Wolke, D.</i> , A note on the Möbius and Liouville functions .....	295
<i>Plonka, J.</i> , On the sum of a $r$ -semilattice ordered system of algebras .....	301
<i>Popenda, J.</i> , On the boundedness of the solutions of second order differential equations .....	89
<i>Poyatos, F.</i> , Archimedean decompositions of left $S$ -semimodules and semirings .....	323
<i>Prodinger, H.</i> and <i>Panny, W.</i> , The expected height of paths for several notions of height .....	119
<i>Puri, M. L.</i> and <i>Chiang, C.-Y.</i> , Tests of subhypotheses in linear regression based on rank-order estimates .....	237
<i>Puri, M. L.</i> and <i>Seoh, M.</i> , Berry—Esséen theorems for signed linear rank statistics under near location alternatives .....	197
<i>Reimnitz, P.</i> , An arcsine-law for the oscillating random walk .....	439
<i>Rosa, A.</i> and <i>Colbourn, C. J.</i> , Indecomposable triple systems with $\lambda=4$ .....	139
<i>Salerno, S.</i> , On the distribution of $x_1^k + \dots + x_n^k$ in the arithmetic progressions .....	357
<i>Salerno, S.</i> and <i>Perelli, A.</i> , On $2k$ -dimensional density estimates .....	345
<i>Seoh, M.</i> and <i>Puri, M. L.</i> , Berry—Esséen theorems for signed linear rank statistics under near location alternatives .....	197
<i>Simon, L.</i> , On approximation of solutions of exterior boundary value problems .....	413
<i>Sitaramachandrarao, R.</i> and <i>Subbarao, M. V.</i> , The distribution of values of a class of arithmetic functions .....	77
<i>Srivastava, K. B.</i> , A remark on Mathur's paper. Simple proof of Telyakovskii—Gopengauz's theorem .....	223
<i>Stein, S.</i> , Lattice-tiling by certain star bodies .....	71
<i>Stern, M.</i> and <i>Damaschke, P.</i> , A characterization of generalized matroid lattices .....	425
<i>Subbarao, M. V.</i> and <i>Sitaramachandrarao, R.</i> , The distribution of values of a class of arithmetic functions .....	77
<i>Veidinger, L.</i> , On the order of convergence of a finite element method for the biharmonic equation .....	255
<i>Wolke, D.</i> , <i>Harman, G.</i> and <i>Pintz, J.</i> , A note on the Möbius and Liouville functions .....	295
<i>Xekalaki, E.</i> , Some bivariate extensions of the generalized Waring distribution .....	173
<i>Xenikakis, P. J.</i> and <i>Miller, H. I.</i> , Some results on the Cantor set .....	309

# ON BROWNIAN BRIDGE AND EXCURSION

J. P. IMHOF

## 1. Introduction

A number of explicit results about extrema over time intervals either fixed or determined by first and last passage times can be obtained easily for Brownian motion or meander and for the three-dimensional Bessel process from joint densities expressed in natural factorized form ([7]). This also provides an elementary direct approach to some path decompositions. This approach is here extended to Brownian bridge and excursion.

First, the relations between these two processes due to Vervaat are established in a simple manner.

Then, starting for Brownian bridge from simultaneous consideration of both maximum and minimum, all known marginals are deduced. Two density factorizations at an extremum first given by Vincze [13] are obtained as particular cases. Some results are mentioned for the time spent above a level.

Brownian excursion is finally considered. Explicit densities upon which Knight [8] based his study of local time are obtained in a simple form equivalent to his via some  $\theta$  function transformation formulas. A density factorization at the maximum is found to transform in this way to a convolution formula first pointed out by Chung [3] and hitherto unexplained.

Processes are considered in the canonical description.  $\Omega$  is the space of continuous functions  $\omega: \mathbf{R}^+ \rightarrow \mathbf{R}$ ,  $\Omega_1$  the one of continuous  $\omega: [0, 1] \rightarrow \mathbf{R}$ . The process is  $X = \{X_t\}$  with  $t \geq 0$ , respectively  $0 \leq t \leq 1$ , given by the variables  $X_t(\omega) = \omega(t)$  which generate the natural  $\sigma$ -fields  $\mathcal{F}_t = \sigma\{X_s: 0 \leq s \leq t\}$ . For  $\Omega$  we set in addition  $\mathcal{F} = \sigma\{X_s: s \geq 0\}$ . On  $(\Omega, \mathcal{F})$ , respectively  $(\Omega_1, \mathcal{F}_1)$ , we call  $P$ , respectively  $P_0$  and  $P_0^+$  the probability laws which make  $X$  be Brownian motion, respectively Brownian bridge and (scaled) Brownian excursion. We further write  $\tau_x = \inf\{s: X_s = x\}$  and

$$m_t = \inf\{X_s: 0 \leq s \leq t\}, \quad M_t = \sup\{X_s: 0 \leq s \leq t\},$$

$$\mu_t = \inf\{s: X_s = m_t\}, \quad \sigma_t = \inf\{s: X_s = M_t\}.$$

For  $t=1$ , these are shortened to  $m, M, \mu, \sigma$ . Formuli are often simpler if one uses for some basic functions a notation giving precedence to functional form over probabilistic meaning. We therefore depart from the notation of [7] and write for all  $t > 0$  and  $x \in \mathbf{R}$ ,  $p_t(x) = (2\pi t)^{-1/2} \exp\{-x^2/2t\}$ ,  $g_t(x) = (2\pi t^3)^{-1/2} x \exp\{-x^2/2t\} = -\frac{\partial}{\partial x} p_t(x)$ . Thus  $p_t(x) = P(X_t \in dx)/dx$  for all  $x$  but  $g_t(x) = P(\tau_x \in dt)/dt$  for

1980 *Mathematics Subject Classification*. Primary 60J65; Secondary 60E05.

*Key words and phrases*. Brownian bridge, Brownian excursion.

$x > 0$  only. Similarly we let for all  $x, y \in \mathbb{R}$ ,  $q_t(x, y) = p_t(x - y) - p_t(x + y)$  so  $q_t(x, y) = P(X_t \in dy - x, \tau_{-x} > t) / dy$  holds for  $xy > 0$  only.

As in [7] the factorization

$$(1.1) \quad P(M_t \in dy, \sigma_t \in ds, X_t \in dx) = 2g_s(y)g_{t-s}(y-x)dydsdx,$$

valid for  $0 < s < t$  and  $0 \vee x < y$ , is basic.

## 2. Relations between Brownian bridge and excursion

The simplest way to obtain the laws of specific functionals of Brownian bridge is often to let  $h \downarrow 0$  in the corresponding ones for Brownian motion over  $0 \leq t \leq 1$ , conditioned to  $X_1 \in [0, h]$  (Billingsley [2]: the continuous mapping theorem applies in all cases we consider). By Scheffé's theorem ([2] p. 224) the passage to the limit can be carried out directly on conditional probability densities if their limit is again a probability density. The following illustrative example indicates the pattern.

For  $0 < s < t < 1$  and  $0 < x < y$ , the Markov property at  $\tau_x$  and (1.1) give

$$\begin{aligned} & P(\tau_x \in ds, \sigma \in dt, M \in dy, X_1 \in [0, h]) / P(X_1 \in [0, h]) = \\ & 2dsdt dy g_s(x)g_{t-s}(y-x) \int_0^h g_{1-t}(y-z)dz / \int_0^h p_1(z)dz. \end{aligned}$$

The limit for  $h \downarrow 0$  of the quotient of integrals is  $g_{1-t}(y)/p_1(0) = \sqrt{2\pi}g_{1-t}(y)$ . Therefore

$$(2.1) \quad P_0(\tau_x \in ds, \sigma \in dt, M \in dy) = 2\sqrt{2\pi}g_s(x)g_{t-s}(y-x)g_{1-t}(y)dsdtdy.$$

Marginalizations show that this is a probability density. For instance,

$$P_0(\sigma \in dt, M \in dy) = 2\sqrt{2\pi}g_t(y)g_{1-t}(y)dt dy.$$

Vincze ([13], Satz 3) had obtained this by passage to the limit from tied down simple symmetric random walk. One may notice the other marginals  $P_0(\tau_x \in ds, M \in dy) = 2\sqrt{2\pi}g_s(x)g_{1-s}(2y-x)dsdy$ ,  $P_0(\tau_x \in ds) = \sqrt{2\pi}g_s(x)p_{1-s}(x)ds$ ,  $P_0(M \in dy) = 2\sqrt{2\pi}g_1(2y)dy$ , this last one at least wellknown.

More generally, the pre-maximum and post-maximum behavior can be detailed by giving a description in terms of finite dimensional events. Let therefore  $0 < r_1 < \dots < r_k < s < t_1 < \dots < t_n < 1$  and define the two formal events

$$C = \{X_{r_i} \in dw_i, i = 1, \dots, k\}, \quad D = \{X_{t_j-s} \in dz_j, j = 1, \dots, n\},$$

so that shift by  $s$  gives  $D \circ \theta_s = \{X_{t_j} \in dz_j, j = 1, \dots, n\}$ . Write also e.g.  $C + a = \{X_{r_i} \in a + dw_i, i = 1, \dots, k\}$ .

LEMMA 1. For  $y$  and all  $w_i, z_j > 0$  there holds

$$(2.2) \quad P_0(C - y, m \in -dy, D \circ \theta_s - y | \mu = s) = P_0^+(D, X_{1-s} \in dy, C \circ \theta_{1-s}).$$

PROOF. Proceeding like for (2.1) and taking into account that the transition from  $X_{t_n} = z_n - y$  to  $X_1 \in [0, h]$  must occur without hit of  $-y$  one obtains after let-



ting  $h \downarrow 0$ ,

$$P_0(C-y, m \in -dy, D \circ \theta_s - y, \mu \in ds) =$$

$$\sqrt{2\pi} P(C-y, m \in -dy, D \circ \theta_s - y, \mu \in ds) q_{1-r_n}(z_n, y).$$

By (1.1) and the Markov property,  $P(\dots)$  above equals

$$2\pi_1 g_{s-r_k}(w_k) g_{t_1-s}(z_1) \pi_2^* dy ds dz_1,$$

where the factors are (with  $r_0=0$  and  $w_0=y$ )

$$\pi_1 = \prod_{i=1}^k q_{r_i-r_{i-1}}(w_{i-1}, w_i) dw_i, \quad \pi_2^* = \prod_{j=1}^{n-1} q_{t_{j+1}-t_j}(z_j, z_{j+1}) dz_{j+1}.$$

As  $\mu$  is uniform the conditioning in the left-hand member of (2.2) amounts to dropping the differential  $ds$ . Rearranging factors this gives

$$P_0(C-y, m \in -dy, D \circ \theta_s - y | \mu = s) = 2\sqrt{2\pi} g_{t_1-s}(z_1) dz_1 \pi_2 \pi_1 g_{s-r_n}(w_n)$$

where  $\pi_2 = \pi_2^* q_{1-t_n}(z_n, y) dy$ . Reference to formula (4.3) of [3] shows this is the right-hand member of (2.2).

The correspondence between Brownian bridge and excursion due to Vervaat is now a simple consequence ([12] for i), oral communication for ii)).

THEOREM 1. i) If  $X$  is Brownian bridge, the process  $U$  defined by

$$U_t = X_{(\mu+t) \pmod{1}} - m, \quad 0 \leq t \leq 1,$$

is Brownian excursion.

ii) If  $\lambda$  is uniform over  $[0, 1]$ , independent of  $X$  which is Brownian excursion, the process  $V$  defined by  $V_t = X_{(\lambda+t) \pmod{1}} - X_\lambda$ ,  $0 \leq t \leq 1$ , is Brownian bridge.

PROOF. i) Integration of (2.2) in  $y$  gives

$$P_0(C-m, D \circ \theta_s - m | \mu = s) = P_0^+(D, C \circ \theta_{1-s}).$$

As the left-hand member equals  $P_0(U_{t_j-s} \in dz_j, j=1, \dots, n, U_{1-s+r_i} \in dw_i, i=1, \dots, k)$ , the  $P_0$ -finite dimensional densities of  $U$  equal the  $P_0^+$ -finite dimensional densities of  $X$ , establishing i).

ii) Because of the assumed independence,

$$(2.3) \quad P_0^+(D, X_{1-s} \in dy, C \circ \theta_{1-s}) = P_0^+(\text{same} | \lambda = 1-s).$$

If  $m^V$  and  $\mu^V$  are the minimum and time of the (first) minimum of  $V$ , the definition of  $V$  shows the above right-hand member equals

$$P_0^+(V_{t_i} \in dw_i - y, i=1, \dots, k, m^V \in -dy, V_{t_j} \in dz_j - y, j=1, \dots, n | \mu^V = s),$$

which together with (2.2) and (2.3) shows that this conditional probability is the same for  $V$  as it is for Brownian bridge. Both  $\mu^V = 1-\lambda$  and the time  $\mu$  for Brownian bridge being uniform, the unconditional probabilities are the same also, i.e. the  $P_0^+$ -finite dimensional distributions of  $V$  equal those of Brownian bridge.

Lemma 1 furthermore shows why the  $P_0$ -density of  $M$  (or equivalently of  $-m$ ) appears as expected occupation time density for Brownian excursion ([3], (6.2) and

[5], (3.4)). If  $S(dy)$  is the time spent by  $X$  in  $dy$  up to time 1, then using (2.2) with  $C$  and  $D$  ignored its  $P_0^+$ -expectation can be written

$$\begin{aligned}\int_0^1 P_0^+(X_t \in dy) dt &= \int_0^1 P_0(-m \in dy | \mu = 1-t) dt = \\ &= \int_0^1 P_0(-m \in dy, \mu \in 1-dt) = P_0(-m \in dy).\end{aligned}$$

### 3. Some results for Brownian bridge

We first need additional notation. For  $x, y, z \in \mathbf{R}$  with  $y$  and  $z \neq 0$  define, all sums being over  $n \in \mathbf{Z}$ ,

$$(3.1) \quad P_t(x, y) = \Sigma p_t(x+2ny), \quad G_t(x, y) = \Sigma g_t(x+2ny),$$

$$(3.2) \quad Q_t(x, y, z) = \Sigma q_t(x, y+2nz) = P_t(x-y, z) - P_t(x+y, z).$$

For  $0 < x < y$  one has the interpretation ([3], Proposition 8)

$$(3.3) \quad G_t(x, y) dt = P(\tau_x \in dt, \tau_{x-y} > t),$$

and for  $0 < x < z, 0 < y < z$  it is wellknown (e.g. [2], (11.10)) that

$$(3.4) \quad Q_t(x, y, z) dy = P^x(X_t \in dy, \tau_0 \wedge \tau_z > t).$$

Here  $P^x$  is the law of Brownian motion starting at  $x$ . It is established in [7] and indicated independently in [10] that (1.1) extends for  $0 < s < t, -z < 0 < y$  and  $-z < x < y$ , to

$$(3.5) \quad P(M_t \in dy, \sigma_t \in ds, X_t \in dx, \tau_{-z} > t) = 2G_s(y, y+z)G_{t-s}(y-x, y+z) dy ds dx.$$

A further basic function is defined for  $t, y > 0$  and arbitrary  $s \in (0, t)$  by

$$(3.6) \quad E_t(y) = \int_0^y G_s(z, y)G_{t-s}(y-z, y) dz.$$

It is shown in [7] that this effectively does not depend on  $s$  and that

$$(3.7) \quad E_t(y) = \frac{\partial}{\partial t} P_t(y, y) = -\frac{1}{2t} \frac{\partial}{\partial y} \{y P_t(y, y)\}.$$

A useful formula as well as a motivation for (3.6) are obtained when considering for  $0 < x < y$ ,  $m_{\tau_x} = \inf \{X_s: 0 \leq s \leq \tau_x\}$  and  $\varrho = \inf \{s < \tau_x: X_s = m_{\tau_x}\}$ . Applying (3.5), the Markov property at a time  $r+s$  where  $r > 0$  and  $0 < s < t-r$ , then (3.3), and finally referring to (3.6), one obtains for Brownian motion

$$P(m_{\tau_x} \in x-dy, \varrho \in dr, \tau_x \in dt) = 2G_r(y-x, y)E_{t-r}(y) dy dr dt.$$

From (3.3) one has furthermore  $P(m_{\tau_x} \in x-dy, \tau_x \in dt) = \frac{\partial}{\partial y} G_t(x, y) dt dy$  and margin-

alization above thus gives

$$(3.8) \quad 2 \int_0^t G_r(y-x, y) E_{t-r}(y) dr = \frac{\partial}{\partial y} G_t(x, y).$$

This can also be checked by term-by-term integration.

Let  $0 < s < t < 1$ ,  $0 < x, y$  and  $-x < z < y$ . Applying the Markov property at some  $u \in (s, t)$  and (3.5) to both the pre and post- $u$  parts of  $X$ , use of (3.6) gives

$$P(m \in -dx, \mu \in ds, M \in dy, \sigma \in dt, X(1) \in dz) = \\ 4G_s(x, x+y) E_{t-s}(x+y) G_{1-t}(y-z, x+y) dx ds dy dt dz.$$

We can now pass to Brownian bridge. Proceeding as with (2.1)

$$(3.9) \quad P_0(m \in -dx, \mu \in ds, M \in dy, \sigma \in dt) = \\ 4\sqrt{2\pi} G_s(x, x+y) E_{t-s}(x+y) G_{1-t}(y, x+y) dx ds dy dt.$$

A first marginalization is accomplished by (3.8):

$$P_0(m \in -dx, M \in dy, \sigma > \mu \in ds) = \\ 2\sqrt{2\pi} G_s(x, x+y) \frac{\partial}{\partial y} G_{1-s}(x, x+y) dx dy ds.$$

Adding the corresponding expression for  $\sigma < \mu$  one has

$$P_0(m \in -dx, \mu \in ds, M \in dy) = 2\sqrt{2\pi} \frac{\partial}{\partial y} \{G_s(x, x+y) G_{1-s}(x, x+y)\} dx ds dy.$$

As  $G_s(x, x) = 0$ , integration in  $y$  gives

$$(3.10) \quad P_0(m \in -dx, \mu \in ds, M < y) = 2\sqrt{2\pi} G_s(x, x+y) G_{1-s}(x, x+y) dx ds.$$

One can notice that marginalization with respect to  $\sigma$  in (3.5) gives, in view of the interpretation (3.4),

$$(3.11) \quad 2 \int_0^t G_s(y, y+z) G_{t-s}(y-x, y+z) ds = \frac{\partial}{\partial y} Q_t(z, x+z, y+z).$$

Using this to integrate the previous equation gives

$$P_0(m > -x, M < y) = \sqrt{2\pi} Q_1(y, y, x+y)$$

which is (11.38) in [2].

Let the maximum absolute deviation over  $[0, 1]$  be  $A = \max \{-m, M\}$ , call  $\alpha$  the time when it (first) occurs. There follows from (3.10), for  $0 < s < 1$  and  $x > 0$

$$P_0(A \in dx, \alpha \in ds) = 4\sqrt{2\pi} G_s(x, 2x) G_{1-s}(x, 2x) dx ds,$$

easily seen to be equivalent to a result of Vincze ([13], Satz 4).

Write furthermore  $R = M - m$ . For  $s < t$ , (3.9) gives

$$P_0(R \in dz, \mu \in ds, \sigma \in dt) = 4\sqrt{2\pi} E_{t-s}(z) \int_0^z G_s(x, z) G_{t-t}(z-x, z) dx \cdot dz ds dt.$$

The integral equals  $E_{1-t+s}(z)$  and for  $s > t$ , one must interchange  $s$  and  $t$  in the right-hand member. If  $\delta = |\sigma - \mu|$  one obtains therefore, for  $0 < u < 1$  and  $z > 0$ ,

$$(3.12) \quad P_0(R \in dz, \delta \in du) = 4\sqrt{2\pi} E_u(z) E_{1-u}(z) dz du.$$

Unfortunately, integration in  $z$  cannot be done term-by-term (a way around the difficulty is indicated for a much simpler case in [5], Corollary (3.2)). The law of  $\delta$  and joint law of  $\mu, \sigma$  (with two uniform marginals) thus remain unknown.

We now mention briefly a few more explicit results for Brownian bridge, ignoring routine calculations. First consider for  $y \in \mathbb{R}$

$$\gamma_y = \sup \{s \leq 1 : X_s = y\}, \quad \Delta_y = \gamma_y - \tau_y,$$

setting  $\Delta_y = 0$  when  $\tau_y > 1$ . Using the Markov property at  $\tau_y$ , then (2.5) of [3] and proceeding as with (2.1) we can write when  $y > 0$ ,

$$P_0(\tau_y \in dt, \gamma_y \in du) = \sqrt{2\pi} g_t(y) p_{u-t}(0) g_{1-u}(y) dt du, \quad 0 < t < u < 1.$$

There results for  $0 < s < 1$

$$P_0(\Delta_y \in ds) = \sqrt{2\pi} p_s(0) g_{1-s}(2y) ds,$$

and use of (2.8) and (3.16) in [3] gives for the  $P_0$ -expectation of  $\Delta_y$  the value  $\sqrt{2\pi} p_1(y)[1 - y\varrho(y)]$ , where  $\varrho(y)$  is Mill's ratio  $\Phi(-y)/p_1(y)$ . ( $\Phi$  is the standard normal distribution function and  $\varrho(y) < 1/y$  is wellknown [11]).

Next, use of the Markov property at  $t$  for Brownian motion, followed by passage to Brownian bridge as with (2.1) gives for  $0 < t < 1$ ,  $0 < y$ ,

$$P_0(M_t < y) = \sqrt{2\pi} \int_{-\infty}^y q_t(y, y-z) p_{1-t}(z) dz.$$

One obtains by elementary integration

$$P_0(\tau_y > t) = P_0(M_t < y) = \Phi\left(\frac{y}{\sqrt{t(1-t)}}\right) - \exp(-2y^2) \Phi\left(\frac{(2t-1)y}{\sqrt{t(1-t)}}\right).$$

#### 4. Some results for Brownian excursion

Theorems 1 and 4 of [3] say that formal time reversal is legitimate for the joint densities considered there. We need another such result and let

$$\gamma_{y,t} = \sup \{s : 0 \leq s \leq t, X_s = y\} \quad (= 0 \text{ if } \tau_y > t).$$

LEMMA 2. For  $0 < x < y$  and  $0 < s < t$ ,

$$(4.1) \quad P(\gamma_{0,t} \in ds, M_t < y, X_t \in dx) = q_s(y, y) G_{t-s}(x, y) ds dx.$$



PROOF. Application of the Markov property at times  $t-s$  and  $s$ , respectively, gives the two equalities

$$\begin{aligned} P^x(\tau_0 > t-s, M_t < y, X_t \in d0)/d0 &= \int_0^y Q_{t-s}(x, z, y) q_s(y-z, y) dz = \\ &= P(\gamma_{0,t} < s, M_t < y, X_t \in dx)/dx. \end{aligned}$$

Taking  $\frac{d}{ds}$  of the first member one obtains a density which, by the strong Markov property at  $\tau_0$ , can be written  $G_{t-s}(x, y) q_s(y, y)$ . This equals  $\frac{d}{ds}$  of the third member, i.e. the left-hand member of (4.1) over  $ds dx$ .

We shall use (4.1) under the equivalent form

$$(4.2) \quad P^x(\gamma_{y,v} \in du, X_v \in dz, m_v > 0) = q_u(y, y) G_{u-v}(y-z, y) du dz,$$

where  $0 < z < y$  and  $0 < u < v$ .

The two joint densities we want to give for Brownian excursion will follow, one from rewriting (3.5) when  $t=1$  in the form

$$(4.3) \quad P^x(M \in dy, \sigma \in ds, X_1 \in dz, m > 0)/dz = 2G_s(y-x, y) G_{1-s}(y-z, y) dy ds,$$

where  $0 < x < y$ ,  $0 < z < y$ ,  $0 < s < 1$ , the other (where  $\gamma_y = \gamma_{y,1}$ ) from

$$(4.4) \quad P^x(\tau_y \in ds, \gamma_y \in dt, X_1 \in dz, m > 0)/dz = G_s(y-x, y) q_{t-s}(y, y) G_{1-t}(y-z, y) ds dt,$$

where  $0 < x < y$ ,  $0 < z < y$  and  $0 < s < t < 1$ . This results from the Markov property at  $\tau_y$  together with (4.2).

It is immediate that  $\left[ \frac{\partial}{\partial x} q_t(x, y) \right]_{x=0} = 2g_t(y)$ . Louchard [10] has observed that comparison of the joint densities for Brownian motion over  $[0, 1]$  starting at  $x > 0$  and subjected to  $m > 0$  with those for Brownian excursion as given in [3] implies, stated here e.g. in the case of  $M, \sigma$ ,

$$(4.5) \quad P_0^+(M \in dy, \sigma \in ds) = \sqrt{\pi/2} \left[ \frac{\partial^2}{\partial x \partial z} P^x(M \in dy, \sigma \in ds, X_1 \in dz, m > 0)/dz \right]_{x=z=0}$$

In order to apply this, we observe that  $\frac{\partial}{\partial x} g_s(y-x+2ny) = 2 \frac{\partial}{\partial s} p_s(y-x+2ny)$  (the heat equation) gives by summation  $\frac{\partial}{\partial x} G_s(y-x, y) = 2 \frac{\partial}{\partial s} P_s(y-x, y)$ , so that according to (3.7)

$$\left( \frac{\partial}{\partial x} G_s(y-x, y) \right)_{x=0} = 2E_s(y).$$

Using this in (4.5) and in the similar relation for  $\tau_y, \gamma_y$  one deduces from (4.3)

and (4.4) the joint densities ( $y > 0$ ),

$$(4.6) \quad P_0^+(M \in dy, \sigma \in ds) = 4\sqrt{2\pi} E_s(y) E_{1-s}(y) dy ds, \quad 0 < s < 1,$$

$$(4.7) \quad P_0^+(\tau_y \in ds, \gamma_y \in dt) = 2\sqrt{2\pi} E_s(y) q_{t-s}(y, y) E_{1-t}(y) ds dt, \quad 0 < s < t < 1.$$

Notice that (4.6) also follows from (3.12) and Theorem 1 i). We state the next result as a theorem.

**THEOREM 2.** For  $y > 0$ ,  $0 < u < 1$ , and  $\Delta_y = \gamma_y - \tau_y$

$$(4.8) \quad P_0^+(\Delta_y \in du)/du = \sqrt{2\pi} q_u(y, y) \frac{\partial^2}{\partial u \partial y} P_{1-u}(0, y).$$

**PROOF.** If one considers the excursion of arbitrary duration  $t$  instead of the scaled excursion and let  $P_0^{+,t}$  be its law, then instead of (4.6) one has for  $y > 0$ ,  $0 < s < t$ ,

$$P_0^{+,t}(M_t \in dy, \sigma_t \in ds) = 4\sqrt{2\pi t^3} E_s(y) E_{t-s}(y) dy ds.$$

Chung [3] gives, in our notation,  $P_0^{+,t}(M_t < y) = -2\sqrt{2\pi t^3} \frac{\partial}{\partial t} P_t(0, y)$ . Marginalization with respect to  $\sigma_t$  above therefore shows that

$$(4.9) \quad 2 \int_0^t E_s(y) E_{t-s}(y) ds = -\frac{\partial^2}{\partial y \partial t} P_t(0, y).$$

Changing to the variables  $\tau_y, \Delta_y$  in (4.7) and marginalizing with respect to  $\tau_y$ , one obtains (4.8).

The general routine for formally writing down joint densities (either for  $P_0$  or for  $P_0^+$ ) factorized at first and last passages as well as extrema is clear from the examples we have considered. For instance,

$$(4.10) \quad P_0^+(\tau_y \in ds) = 2\sqrt{2\pi} E_s(y) g_{1-s}(y) ds, \quad 0 < y, \quad 0 < s < 1.$$

Furthermore, insertion of finite dimensional events as in Lemma 1 easily permits to identify pieces of processes, only the writing becoming cumbersome. Here is one example.

**THEOREM 3.** Given  $\Delta_y = t$  for scaled Brownian excursion  $X$  the process  $U_s = X_{\tau_y+s}$ ,  $0 \leq s \leq t$  is a Brownian bridge of duration  $t$  conditioned to a minimum  $> -y$ .

Using a completely different approach, Knight [8] has obtained for the densities in (4.8) and (4.10) expressions not obviously equivalent to ours. We establish this equivalence and notice an interesting consequence.

A classical transformation formula for theta functions is ([1], (17.1)), all sums below being over  $n \in \mathbb{Z}$ :

$$\sum e^{-\alpha n^2} = \sqrt{\pi/\alpha} \sum e^{-\pi^2 n^2/\alpha}, \quad \alpha > 0.$$

Using this twice in  $\sum (-1)^n \exp(-\alpha n^2) = 2\sum \exp(-4\alpha n^2) - \sum \exp(-\alpha n^2)$ , the result

can be written

$$\Sigma e^{-\beta(2n+1)^2} = \frac{1}{2} \sqrt{\pi/\beta} \Sigma (-1)^{-n} e^{-n^2\pi^2/4\beta}, \quad \beta > 0.$$

The functions  $P_t(0, y)$  and  $P_t(y, y)$  therefore have the alternative expressions

$$(4.11) \quad P_t(0, y) = \frac{1}{2y} \Sigma e^{-n^2\pi^2 t/2y^2}, \quad P_t(y, y) = \frac{1}{2y} \Sigma (-1)^n e^{-n^2\pi^2 t/2y^2}.$$

In terms of the distribution functions (for  $x > 0$ )

$$F_1(x) = \Sigma (-1)^n e^{-n^2 x}, \quad F_2(x) = \Sigma (1 - 2n^2 x) e^{-n^2 x},$$

and their densities  $f_1, f_2$ , one has consequently for  $y, t > 0$ ,

$$(4.12) \quad -\frac{\partial}{\partial y} P_t(0, y) = \frac{1}{2y^2} F_2\left(\frac{\pi^2 t}{2y^2}\right), \quad -\frac{\partial^2}{\partial t \partial y} P_t(0, y) = \frac{\pi^2}{4y^4} f_2\left(\frac{\pi^2 t}{2y^2}\right),$$

$$(4.13) \quad P_t(y, y) = \frac{1}{2y} F_1\left(\frac{\pi^2 t}{2y^2}\right), \quad E_t(y) = \frac{\partial}{\partial t} P_t(y, y) = \frac{\pi^2}{4y^3} f_1\left(\frac{\pi^2 t}{2y^2}\right).$$

With this, our densities (4.10) and (4.8) are easily seen to be equal respectively to 1/2 and 1/4 of the densities in Theorem 1.1 and Corollary 1.3 of [8].

Those correction factors have been pointed out by Knight [9]. Other interesting considerations about  $F_1$  are in [4].

If one substitutes in (4.9) according to (4.12) and (4.13) and sets  $y = \pi/\sqrt{2}$  there comes

$$\int_0^t f_1(s) f_1(t-s) ds = f_2(t).$$

Chung [3] first observed that  $f_2 = f_1 * f_1$  and Knight [8] has given a proof by characteristic functions. The above provides a probabilistic explanation for this convolution.

\* We conclude by pointing out a further connexion. One has from (3.11) for all  $t > 0$ , with  $y > 0$ ,

$$4 \int_0^t G_s(y, 2y) G_{t-s}(y, 2y) ds = 2 \left[ \frac{\partial}{\partial y} Q_t(z, z, y+z) \right]_{z=y} = \frac{\partial}{\partial y} Q_t(y, y, 2y),$$

where the second equality is easily checked directly. On the other hand (3.3) implies

$$P(\inf \{s > 0: |X_s| = y\} \in dt) = 2G_t(y, 2y) dt,$$

and the corresponding L. T. (Laplace transform) is known to be  $1/\cosh(y\sqrt{2\lambda})$ ,  $\lambda > 0$ . The above convolution therefore shows that for  $y > 0$ ,

$$f(t) = \frac{\partial}{\partial y} Q_t(y, y, 2y) = (2\pi t)^{-1/2} \frac{\partial}{\partial y} \Sigma (-1)^n e^{-n^2\pi^2 y^2/t}$$

is a probability density over  $(0, \infty)$  with L. T.  $(\cos y\sqrt{2\lambda})^{-2}$ . For  $y=1$  this is the transform obtained in [6] for a limit law concerning Brownian excursion occupation time. We shall show elsewhere that  $f(t)$  is also the density for the hitting time of  $2y$  by the range of Brownian motion.

## REFERENCES

- [1] BELLMAN, R., *A brief introduction to theta functions*. Holt, Rinehart and Winston, New York, 1961. *MR* 23 # A 2556.
- [2] BILLINGSLEY, P. *Convergence of Probability Measures*, John Wiley & Sons, Inc., New York—London—Sydney, 1968. *MR* 38 # 1718.
- [3] CHUNG, K. L., Excursions in Brownian motion, *Ark. Mat.* 14 (1976), 155—177. *MR* 57 # 7791.
- [4] CHUNG, K. L., A cluster of great formulas, *Acta Math. Acad. Sci. Hungar.* 39 (1982), 65—67. *MR* 83g: 60029.
- [5] DURRETT, R. T. and IGLEHART, D. L., Functionals of Brownian meander and Brownian excursion, *Ann. Probability* 5 (1977), 130—135. *MR* 55 # 9301.
- [6] GETOOR, R. K. and SHARPE, M. J., Excursions of Brownian motion and Bessel processes, *Z. Wahrsch. Verw. Gebiete* 47 (1979), 83—106. *MR* 80b: 60104.
- [7] IMHOF, J. P., Density factorizations for Brownian motion, meander and the three-dimensional Bessel process, and applications, *J. Appl. Prob.* 21 (1984), 500—510.
- [8] KNIGHT, F., On the excursion process of Brownian motion, *Trans. Amer. Math. Soc.* 258 (1980), 77—86. *MR* 81d: 60081.
- [9] KNIGHT, F., On the excursion process of Brownian motion, *Zentralblatt für Math.* 426 (1980), 60073 (Autorreferat for [8]).
- [10] LOUCHARD, G., Kac's formula, Levy's local time and Brownian excursion, *J. Appl. Prob.* 21 (1984), 479—499.
- [11] SHENTON, L. R., Inequalities for the normal integral including a new continued fraction, *Biometrika*, 41 (1954), 177—189. *MR* 15—884.
- [12] VERVAAT, W., A relation between Brownian bridge and Brownian excursion, *Ann. Probab.* 7 (1979), 143—149. *MR* 80b: 60107.
- [13] VINCZE, I., Einige zweidimensionale Verteilungs- und Grenzverteilungssätze in der Theorie der geordneten Stichproben. *Magyar Tud. Akad. Mat. Kutató Int. Közl.* 2 (1957), 183—209. *MR* 21 # 3915.

(Received February 28, 1983)

SECTION DE MATHÉMATIQUES  
UNIVERSITÉ DE GENÈVE  
CASE POSTALE 240  
CH—1211 GENÈVE 24  
SWITZERLAND



# AN ASYMPTOTIC STATEMENT CONCERNING THE SOLUTIONS OF THE DIFFERENTIAL EQUATION $x'' + a(t)x = 0$

AN ALTERNATIVE PROOF OF A THEOREM BY G. PRODI AND G. TREVISAN

I. BIHARI

If in the differential equation

$$(1) \quad x'' + a(t)x = 0, \quad t \in I = [t_0, \infty), \quad t_0 \in \mathbb{R}, \quad ' = \frac{d}{dt}$$

function  $a(t)$  has the following properties:  $a(t)$  is increasing and

$$a(t) > 0, \quad a(t) \in C_1(I), \quad \lim_{t \rightarrow \infty} a(t) = \infty$$

then — corresponding to the Prodi—Trevisan theorem [1—2] in question — equation (1) admits at least one non-trivial solution  $x(t)$  with  $\lim_{t \rightarrow \infty} x(t) = 0$ . It is clear that this  $x(t)$  is the unique such solution — disregarded its multiples — provided there is at all a solution not tending to zero. — We give here a *new proof of this theorem*.

PROOF. Suppose the contrary and let  $x_i(t)$  ( $i=1, 2$ ) be a pair of linear independent solutions of (1). For an arbitrary solution  $x(t)$  of (1) define the function  $A(t)$  by

$$(2) \quad A(t) = x(t)^2 + \frac{x'(t)^2}{a(t)}.$$

Since  $A' = -\frac{a'}{a^2} x'^2 \leq 0$ ,  $A(t)$  decreases and  $\lim_{t \rightarrow \infty} A(t) = A > 0$ . Let  $A_i(t)$  and  $A_i$  be defined by

$$A_i(t) = x_i(t)^2 + \frac{x_i'(t)^2}{a(t)}, \quad \lim_{t \rightarrow \infty} A_i(t) = A_i > 0, \quad i = 1, 2.$$

Furthermore let  $x(t)$  be

$$(3) \quad x(t) = c_1 x_1(t) + c_2 x_2(t), \quad (c_i = \text{const}; i = 1, 2).$$

We shall show that the values of  $c_i$  ( $i=1, 2$ ) can be determined in such a way that  $\lim_{t \rightarrow \infty} x(t) = 0$  which, of course, will contradict to our starting assumption. Substitute (3) in (2). Then with the notation

$$A_{12}(t) = x_1(t) \cdot x_2(t) + \frac{x_1'(t)x_2'(t)}{a(t)}$$

we have

$$(4) \quad A(t) = A_1(t)c_1^2 + 2A_{12}(t)c_1c_2 + A_2(t)c_2^2.$$

Since the limits

$$\lim_{t \rightarrow \infty} A_i(t) = A_i \quad (i = 1, 2), \quad \lim_{t \rightarrow \infty} A(t) = A$$

exist (and are positive), therefore by (4)  $\lim_{t \rightarrow \infty} A_{12}(t) = A_{12}$  exists, too, and by (4) the relation

$$(5) \quad A = A_1 c_1^2 + 2A_{12} c_1 c_2 + A_2 c_2^2$$

holds. The determinant of this quadratic form is

$$D = \begin{vmatrix} A_1 & A_{12} \\ A_{12} & A_2 \end{vmatrix} = A_1 A_2 - A_{12}^2$$

while the determinant of the quadratic form in (4) is

$$D(t) = A_1(t) \cdot A_2(t) - A_{12}^2(t) = \frac{(x_1' x_2 - x_2' x_1)^2}{a(t)} = \frac{k^2}{a(t)} > 0 \quad (k = \text{const})$$

being the Wronskian of  $x_1(t)$ ,  $x_2(t)$  constant.<sup>1</sup> Consequently, there is a unique ratio

$$c_1 : c_2 = -\frac{A_{12}}{A_1} = -\frac{A_2}{A_{12}},$$

for which  $A=0$  and  $\lim_{t \rightarrow \infty} x(t)=0$  where, of course

$$(6) \quad x(t) = \lambda(A_{12}x_1(t) - A_1x_2(t)), \quad \lambda \in \mathbb{R}.$$

Herewith we came to a contradiction. — Another solution, linear independent from  $x(t)$  cannot exist except every solution behaves in the same way.

REMARK 1. Also in the case where  $x_1(t) \rightarrow 0$ ,  $x_2(t) \rightarrow 0$  as  $t \rightarrow \infty$ , the solution  $x(t)$  given by (6) tends to zero as  $t \rightarrow \infty$ . Namely, if the maximum-points of  $|x_1(t)|$  are  $t_k$  ( $k=1, 2, \dots$ ) then we have

$$A_1 = \lim_{t \rightarrow \infty} A_1(t) = \lim_{k \rightarrow \infty} A_1(t_k) = \lim_{k \rightarrow \infty} x_1^2(t_k) = 0.^2$$

REMARK 2. The condition that  $a(t)$  is monotonic can be replaced by the assumption that  $\lim_{t \rightarrow \infty} A(t)$  exists for every solution  $x(t)$  of (1). Then by (4)  $\lim_{t \rightarrow \infty} A_{12}(t) = A_{12}$  exists, too, and the proof goes along the same lines. It would be desirable to find a criterion concerning immediately  $a(t)$  which would involve this weaker assumption.

REMARK 3. It is also well known that every solution of (1) tends to zero as  $t \rightarrow \infty$  provided  $\log a(t)$  tends to infinity "regularly" as  $t \rightarrow \infty$  and  $a(t)$  is monotonic [3]. Moreover this theorem has several extensions to linear and nonlinear second order differential equations [4—5—6].

<sup>1</sup> Thus  $D = \lim_{t \rightarrow \infty} D(t) = 0$  and the quadratic form in (5) is semidefinite.

<sup>2</sup>  $A_1 = 0$  involves  $\frac{x_1'^2}{a} = o(1)$  or  $x_1' = o(\sqrt{a})$ ,  $t \rightarrow \infty$ , while for an  $x(t)$  with  $A > 0$  we have  $\frac{x'^2}{a} < A + \varepsilon$ ,  $|x'(t)| < \sqrt{(A + \varepsilon)a(t)}$ ,  $\varepsilon > 0$ ,  $t \rightarrow \infty$ .

## REFERENCES

- [1] PRODI, G., Un'osservazione sugli integrali dell'equazione  $y'' + A(x)y = 0$  nel caso  $A(x) \rightarrow \infty$  per  $x \rightarrow \infty$ , *Atti Accad. Naz. Lincei. Rend. Cl. Sci. Fis. Mat. Nat.* (8) 8, (1950), 462—464. *MR* 12—334, 1002.
- [2] TREVISAN, G., Sull'equazione differenziale  $y'' + A(x)y = 0$ , *Rend. Sem. Mat. Univ. Padova* 23 (1954), 340—342. *MR* 16—589.
- [3] SANSONE, G., *Equazione differenziali nel campo reale*, Vol. 2, 2nd ed. Nicola Zanichelli, Bologna, 1949. *MR* 11—32.
- [4] OPIAL, Z., Sur l'équation différentielle  $u'' + a(t)u = 0$ , *Ann. Polon. Math.* 5 (1958), 77—93. *MR* 20 # 4047.
- [5] BIHARI, I., Extension of a theorem of Armellini—Tonelli—Sansone to the nonlinear equation  $u'' + a(t)f(u) = 0$ . *Magyar Tud. Akad. Mat. Kutató Int. Közl.* 7 (1962), 63—68. *MR* 26 # 6498.
- [6] BIHARI, I., Asymptotic result concerning equation  $x''|x'|^{n-1} + a(t)x^n = 0$ . Extension of a theorem by Armellini—Tonelli—Sansone, *Studia Sci. Math. Hungar.* 19 (1984), 151—157.

(Received April 12, 1983)

MTA MATEMATIKAI KUTATÓ INTÉZETE  
P.O. BOX 127  
H-1364 BUDAPEST  
HUNGARY





# NOTE TO AN EXTENSION OF A STURMIAN COMPARISON THEOREM

I. BIHARI

This note aims at completing an earlier paper [1] of the author in which beside a Sturmian theorem, some oscillatory and monotonicity properties of the equation

$$y'' + p(x)f(y)f(y') = 0, \quad y = y(x), \quad ' = \frac{d}{dx}$$

were stated under appropriate conditions concerning the functions  $p(x)$ ,  $f(y)$ ,  $g(y')$ .

Denote the class of positive continuous functions  $F(u)$  ( $u \in \mathbf{R}$ ) non decreasing for  $u < 0$  and non increasing for  $u > 0$  by  $\mathcal{F}$ . Consider the pair of second order ordinary nonlinear differential equations

$$y_i'' + p_i(x)f(y_i)g(y_i') = 0, \quad x \in I = [x_0, \infty), \quad x_0 \in \mathbf{R},$$

(1)

$$y_i = y_i(x), \quad ' = \frac{d}{dx}, \quad i = 1, 2$$

under the conditions

(i)

$$p_i \in C(I), \quad p_i > 0, \quad p_1 > p_2, \quad x \in I,$$

(ii)

$$f \in C(\mathbf{R}), f' \neq 0, \quad f(0) = 0, \quad f \in \text{Lip}(1), \quad \frac{f(y)}{y} = O(1)(y \rightarrow 0), \quad \frac{f(y)}{y} \in \mathcal{F}, \quad y \in \mathbf{R}.$$

(iii)

$$g \in \text{Lip}(1), \quad g \in \mathcal{F}.$$

**THEOREM.** Let  $y_i(x)$  ( $i=1, 2$ ) be solutions of the equations (1) for  $x \geq a \geq x_0$ , respectively, satisfying the conditions

$$y_i(a) = A \geq 0 \quad (i = 1, 2), \quad y_1(b) = 0, \quad y_1 \neq 0, \quad a < x < b, \quad x_0 \leq a < b$$

$$(\text{or } y_1(x) > 0 \text{ for } x > a), \quad y_2'(a) \geq y_1'(a) \geq 0.$$

Then

1°

$$x_1 < x_2$$

where  $x_1$  means the abscissa of the maximum point of  $y_1$  next to  $a$  and  $x_2$  — if any — has the same meaning concerning  $y_2$ ,

2°

$$\frac{y_2}{y_1} \nearrow, \quad y_2 > y_1, \quad a < x \leq b \quad (\text{or } x > a).$$

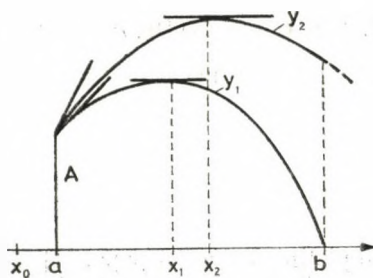


Fig. 1

PROOF. Three cases must be distinguished.

- (I)  $y_2'(a) > y_1'(a)$ ,  
 (II)  $y_2'(a) = y_1'(a) > 0$  (or  $< 0$ )  
 (III)  $y_2'(a) = y_1'(a) = 0$ .

First the Cases (I)–(II) will be treated. It will be shown that in both cases

$$\left. \begin{aligned} y_2'(x) &> y_1'(x) \\ y_2(x) &> y_1(x) \end{aligned} \right\} \text{ for } x > a, x-a \text{ small enough.}$$

This statement is obvious in the first case. In the second one applying the finite Taylor formula to  $y_i(x)$

$$y_i(x) = A + y_i'(a)(x-a) + \frac{1}{2} y_i''(a + \theta_i(x-a))(x-a)^2, \quad 0 < \theta_i < 1$$

$$= A + y_i'(a)(x-a) + \begin{cases} O((x-a)^2), & A \neq 0 \\ O((x-a)^3), & A = 0 \end{cases} \quad (x \rightarrow a+0),$$

(i.e. if  $A=0$  then by (I)  $y_i'' \approx y_i \approx (x-a)$  giving

$$y_2(x) - y_1(x) = \begin{cases} O((x-a)^2), & A \neq 0 \\ O((x-a)^3), & A = 0 \end{cases} \quad (x \rightarrow a+0).$$

In the same way

$$y_i'(x) = y_i'(a) + y_i''(a + \bar{\theta}_i(x-a))(x-a), \quad 0 < \bar{\theta}_i < 1$$

$$y_2'(x) - y_1'(x) = \begin{cases} O((x-a)), \\ O((x-a)^2). \end{cases}$$

Since  $f$  and  $g$  are Lipschitzian we have

$$|f(y_2) - f(y_1)| \leq K_1 |y_1 - y_2| = \begin{cases} O((x-a)^2) \\ O((x-a)^3) \end{cases} \Rightarrow f(y_2) = f(y_1) + \begin{cases} O((x-a)^2) \\ O((x-a)^3) \end{cases},$$

$$|g(y_2') - g(y_1')| \leq K_2 |y_1' - y_2'| = \begin{cases} O((x-a)), \\ O((x-a)^2) \end{cases} \Rightarrow g(y_2') = g(y_1') + \begin{cases} O((x-a)), \\ O((x-a)^2) \end{cases}.$$

whence

$$f(y_2)g(y_2') = f(y_1)g(y_1') + \begin{cases} O((x-a)), & A \neq 0, \\ O((x-a)^3), & A = 0. \end{cases} \quad (x \rightarrow a+0)$$

Therefore by (1)

$$\begin{aligned} y_2'' - y_1'' &= p_1 f(y_1)g(y_1') - p_2 f(y_2)g(y_2') = \\ &= (p_1 - p_2)f(y_1)g(y_1') + \begin{cases} O((x-a)), & A \neq 0, \\ O((x-a)^3), & A = 0, \end{cases} \quad (x \rightarrow a+0) \end{aligned}$$

involving

$$(2) \quad y_2'' > y_1'', \quad y_2' > y_1', \quad y_2 > y_1 > 0$$

for  $x-a > 0$  small enough. Making use again of (1) we have for  $\Delta(x) = y_2'y_1 - y_1'y_2$

$$\Delta' + p_2 f(y_2)g(y_2')y_1 - p_1 f(y_1)g(y_1')y_2 = 0,$$

whence

$$(3) \quad \Delta(x) = \int_a^x \mathcal{J}(x) dx,$$

where

$$\mathcal{J}(x) = \left[ p_1 \frac{f(y_1)}{y_1} g(y_1') - p_2 \frac{f(y_2)}{y_2} g(y_2') \right] y_1 y_2.$$

By virtue of (2) and conditions (i)—(iii)  $\mathcal{J}(x)$  is positive for  $x-a > 0$  and small enough (i.e. in a right-hand neighbourhood of  $a$ ). Thus  $\Delta(x) > 0$  in the same neighbourhood. We state

$$\Delta(x) > 0, \quad a < x < b.$$

In the opposite case there would exist a first (smallest)  $c > a$  with  $\Delta(c) = 0$ . It will be shown that such a number  $c$  does not exist. Suppose that it exists. Then

$$\Delta(x) > 0 \Rightarrow \left( \frac{y_2}{y_1} \right)' = \frac{\Delta(x)}{y_1^2} > 0 \Rightarrow \frac{y_2}{y_1} \nearrow, \quad a < x < c.$$

However,

$$\frac{y_2}{y_1} \Big|_{x=a} = 1 \quad \text{or} \quad \lim_{x \rightarrow a+0} \frac{y_2}{y_1} = \frac{y_2'(a)}{y_1'(a)} = 1,$$

thus  $\frac{y_2}{y_1} > 1$  and  $y_2 > y_1 > 0 (a < x \leq c)$ . Thus

$$(4) \quad \frac{y_2'}{y_2} > \frac{y_1'}{y_1} \quad (a < x < c)$$

and  $y_2' > y_1' > 0$  for  $x-a > 0$  and small enough. This either remains valid up to  $x=c$  (i.e.  $c < x_i, i=1, 2$ ) or  $y_2' > 0 > y_1'$  holds already somewhere before  $c$  involving  $x_1 < x_2$ . The first case leads by (3) to  $\Delta(c) > 0$ . In the second one either  $x_1 < c < x_2$  which is excluded by

$$\frac{y_1'(c)}{y_1(c)} = \frac{y_2'(c)}{y_2(c)} \Rightarrow \operatorname{sgn} y_1'(c) = \operatorname{sgn} y_2'(c)$$

or  $x_2 < c < b$ . Since  $\frac{y_2}{y_1}$  increased up to  $c$

$$y_2(c) > y_1(c) \Rightarrow |y_2'(c)| > |y_1'(c)|,$$

moreover

$$\operatorname{sgn} y_2'(c) = \operatorname{sgn} y_1'(c) < 0,$$

thus we have  $\Delta'(c) = \mathcal{J}(c) > 0 \Rightarrow \Delta'(x) > 0$  for  $c - x > 0$  and small enough, while  $\Delta(c) = 0$  requires  $\Delta'(x) \leq 0$  in the same neighbourhood.

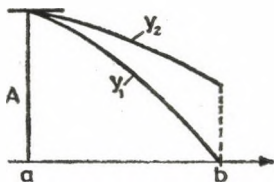


Fig. 2

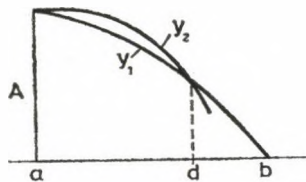


Fig. 3

Case (III) Now

$$(5) \quad \begin{aligned} y_1(a) &= y_2(a) = A \neq 0, \\ y_1'(a) &= y_2'(a) = 0. \end{aligned}$$

Assume e.g.  $A > 0$ , then

$$y_i(x) = A + \frac{1}{2} y_i''(a + \theta_i(x-a))(x-a)^2 = A + O((x-a)^2),$$

$$y_i'(x) = y_i''(a + \bar{\theta}_i(x-a))(x-a) = O((x-a)),$$

$$y_1(x) - y_2(x) = O(x-a)^2 \Rightarrow f(y_1) = f(y_2) + O((x-a)^2),$$

$$y_1' - y_2' = O((x-a)) \Rightarrow g(y_1') = g(y_2') + O((x-a)),$$

$$f(y_1)g(y_1') = f(y_2)g(y_2') + O((x-a)),$$

$$y_2'' - y_1'' = (p_1 - p_2)f(y_1)g(y_1') + O((x-a)), \quad x \rightarrow a+0$$

giving

$$0 > y_2'' > y_1'', \quad 0 > y_2' > y_1',$$

$$|y_1''| > |y_2''|, \quad |y_1'| > |y_2'|, \quad y_2 > y_1 > 0$$

for  $x-a > 0$  small enough, and in the present case the above argument connected with (3) cannot be carried out, however the result remains valid, now too. Namely, if the theorem were false, then a first point  $x=d$  ( $a < d < b$ ) would exist with

$$(6) \quad D = y_2(d) = y_1(d), \quad y_2'(d) \leq y_1'(d) < 0.$$

In order to apply the part already proved of the Theorem (the result in Cases (I) and (II)) reflect Figure 3 to the ordinate  $x=d$ , i.e. introduce the variable  $\xi$  by  $x=d-\xi$ ,

$0 \leq \xi \leq d-a$ , then

$$y_i(x) = z_i(\xi), \quad y'_i = -z'_i, \quad y''_i = z''_i, \quad p_i(x) = q_i(\xi), \quad q_1 > q_2,$$

$$(1') \quad z''_i + q_i(\xi)f(z_i)g(-z'_i) = 0, \quad i = 1, 2$$

and (6) turns to

$$(6') \quad z_1(0) = z_2(0) = D,$$

$$z'_2(0) \equiv z'_1(0) > 0.$$

Since all the conditions of Case (I) and (II) are satisfied concerning (1'), (6') and  $\xi=0$ , we have

$$z_2(\xi) > z_1(\xi), \quad 0 < \xi \leq d-a,$$

whence

$$z_2(d-a) > z_1(d-a) \quad \text{or} \quad y_2(a) > y_1(a)$$

contradicting (5).

The Case  $y_1(a)=y_2(a)$ ,  $y'_2(a) \leq y'_1(a) < 0$  can be settled by the same argument. The Case  $y_1(a)=y_2(a)$ ,  $y'_2(a) > 0$ ,  $y'_1(a) < 0$  can be composed of two of the above cases by inserting a third solution  $y_3(x)$  under the conditions  $y_3(a)=A$ ,  $y'_3(a)=0$ . Comparing  $y_2$  with  $y_3$  and  $y_3$  with  $y_1$ , respectively, we have

$$\{y_2 > y_3, y_3 > y_1\} \Rightarrow y_2 > y_1 \quad (a < x \leq b).$$

REMARK. If  $p_1=p_2$  (the case of a unique equation) and  $y'_2(a) > y'_1(a)$  the result remains valid giving rise to the configurations of Figure 4 and 5 not existing in the case of a linear equation.

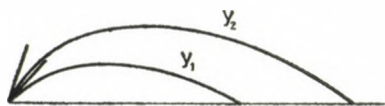


Fig. 4

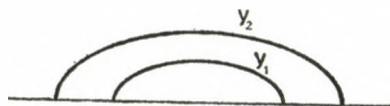


Fig. 5

## REFERENCE

- [1] BIHARI, I., Oscillation and monotony theorems concerning nonlinear differential equations of the second order, *Acta Math. Acad. Sci. Hungar.* 9 (1958), 83—104. MR 20 # 1824.

(Received April 12, 1983)

MTA MATEMATIKAI KUTATÓ INTÉZETE  
P. O. BOX 127  
H-1364 BUDAPEST  
HUNGARY





## ON COMMUTATIVITY OF QUASI-COMMUTATIVE RINGS

H. E. BELL<sup>1</sup>

Let  $Z[X]$  denote the ring of polynomials in one indeterminate with integer coefficients, and let  $P$  be a non-empty subset of  $Z[X]$ . A ring  $R$  will be called *P-commutative* if for each  $x, y \in R$ , there exists a corresponding  $f(X) \in P$  for which  $xy = f(yx)$ , it being understood that  $f(X)$  must have zero constant term if  $R$  does not have  $1$ . If  $R$  is  $Z[X]$ -commutative, it will be called *quasi-commutative*. The first aim of this paper is to study the commutativity behavior of arbitrary quasi-commutative rings (usually with  $1$ ); and we find that in any quasi-commutative  $R$  with  $1$ , commutators must commute. We then proceed to identify certain  $P$  for which all  $P$ -commutative rings with  $1$  must be commutative; and as a corollary, we obtain the result that a ring  $R$  with  $1$  must be commutative if for each  $x, y \in R$ , there exist an integer  $k = k(x, y)$  and a non-negative integer  $n = n(x, y)$  for which  $xy = k(yx)^n$ .

Throughout the paper,  $Z$  will denote the ring of integers. For arbitrary rings  $R$ ,  $\mathcal{C}(R)$  will denote the commutator ideal,  $C$  the center, and  $R^+$  the additive group. The commutator  $xy - yx$  will be denoted by  $[x, y]$ ; and if  $R$  has  $1$ , for  $x \in R$ , the symbol  $x^0$  will be understood as  $1$ . By the *left quasi-center*  $C_l = C_l(R)$  we shall mean the set of all  $y \in R$  with the property that for each  $x \in R$ , there exists  $n \in Z$  for which  $yx = nx$ ; the *right quasi-center*  $C_r = C_r(R)$  is defined analogously. Finally, if  $f(X) = a_0 + a_1X + \dots + a_nX^n \in Z[X]$ , the smallest  $k$  for which  $a_k \neq 0$  will be called the *co-degree* of  $f(X)$ .

## 1. Preliminaries

We begin by stating some necessary results from previous work.

LEMMA 1 [1, Theorem 1]. *Let  $R$  be a ring satisfying an identity  $q(X) = 0$ , where  $q(X)$  is a polynomial in a finite number of non-commuting indeterminates, its coefficients being integers with highest common factor 1. If there exists no prime  $p$  for which the ring of  $2 \times 2$  matrices over  $GF(p)$  satisfies  $q(X) = 0$ , then  $\mathcal{C}(R)$  is nil and the nilpotent elements form an ideal.*

LEMMA 2 [4, Theorem 2]. *If  $R$  has  $1$ , then  $C_l = C_r = C$ .*

Clearly all quasi-commutative rings satisfy the identity  $[xy, yx] = 0$ ; thus, an easy application of Lemma 1 gives

<sup>1</sup> Supported by the Natural Sciences and Engineering Research Council of Canada, Grant No. A 3961.

1980 *Mathematics Subject Classification*. Primary 16A70.

*Key words and phrases*. Commutativity theorems, quasi-commutative rings.

**THEOREM 1.** *If  $R$  is any quasi-commutative ring,  $\mathcal{C}(R)$  is nil and the nilpotent elements form an ideal.*

Of course, the existence of non-commutative nil rings satisfying the identity  $xy + yx = 0$  shows that quasi-commutativity does not in general imply commutativity. Even in rings with 1, quasi-commutativity does not guarantee commutativity, as the following example makes clear.

**EXAMPLE 1.** Let  $p$  be an arbitrary prime. Let  $R^+$  have basis  $\{1, u, v\}$  with  $p^2 \cdot 1 = pu = pv = 0$ ; and define multiplication by  $lu = ul = u$ ,  $lv = vl = v$ ,  $l^2 = 1$ ,  $u^2 = uv = v^2 = 0$ ,  $vu = p \cdot 1$ . Then for each  $x, y \in R$ , there exists a polynomial  $f(X)$  of form  $k + X \in \mathbb{Z}[X]$  for which  $xy = f(yx)$ ; but  $R$  is clearly not commutative. Incidentally, this example shows that in Theorem 7, the definition of  $P_2$  is more natural than it might appear to be.

We need one further lemma.

**LEMMA 3.** *Let  $R$  be an arbitrary ring with 1. Then*

- (a)  *$R$  satisfies the identity  $[xy, yx] = 0$  if and only if it satisfies the identity  $[[x, y], y] = 0$ ;*  
 (b) *if  $R$  satisfies the identity  $[xy, yx] = 0$ , it must also satisfy the identity  $[x, y]^2 = 0$ .*

**PROOF.** (a) If  $R$  satisfies the identity  $[xy, yx] = 0$ , then for each  $x, y \in R$  we have  $0 = [(x + 1)y, y(x + 1)] = [xy + y, yx + y] = [xy, y] + [y, yx] = [x, y]y + y[y, x] = [x, y]y - y[x, y]$ . Conversely, if  $R$  satisfies the identity  $[[x, y], y] = 0$ , it follows that  $[xy, yx] = x[y, yx] + [x, yx]y = xy[y, x] + [x, y]xy = -xy[x, y] + xy[x, y] = 0$ .

(b) We shall use the well-known fact that in rings  $R$  satisfying the identity  $[[x, y], y] = 0$ ,  $[x, y^n] = ny^{n-1}[x, y]$  for all  $x, y \in R$  and all positive integers  $n$ . Thus  $2xy[x, y] = x(2y)[x, y] = x[x, y^2] = [x, y^2]x = 2y[x, y]x = 2yx[x, y]$ , whence

$$(1) \quad 2[x, y]^2 = 0.$$

Similarly,  $3y[x, y]yx = 3y^2[x, y]x = [x, y^3]x = x[x, y^3] = 3xy^2[x, y] = 3xy[x, y]y = 3xy[xy, y] = 3[xy, y]xy = 3[x, y]yxy = 3y[x, y]xy$ ; hence  $3y[x, y]^2 = 0$  for all  $x, y \in R$ . Replacing  $y$  by  $y + 1$  in this last equation and then subtracting the original yields  $3[x, y]^2 = 0$ , which together with (1), forces  $[x, y]^2 = 0$ .

## 2. Commutativity results for special quasi-commutative rings

Before proceeding to our major results, we present some uncomplicated commutativity theorems for some restricted quasi-commutative rings.

**THEOREM 2.** *Let  $f(X) \in \mathbb{Z}[X]$ , and let  $R$  be any ring with 1 which satisfies the identity  $xy = f(yx)$ . Then  $R$  is commutative.*

**PROOF.** Letting  $f_1(X) = f(X) - X$ , we see that  $R$  satisfies the identity  $xy = yx + f_1(yx)$ . Substituting  $a$  for  $x$  and 1 for  $y$  shows that  $f_1(a) = 0$  for all  $a \in R$ , hence  $R$  satisfies the identity  $xy = yx$ .

**THEOREM 3.** *Let  $P_0$  denote the set of all polynomials in  $\mathbb{Z}[X]$  of co-degree at least 1. Then every  $P_0$ -commutative ring with 1 satisfies the identity  $[x, y], w = 0$ .*

PROOF. Let  $x, y$ , and  $w$  be arbitrary elements of  $R$ . There exists  $f(X) \in P_0$  such that  $[x, y]w[x, y] = f(w[x, y]^2)$ ; and since  $[x, y]^2 = 0$  by Lemma 3(b), we have  $[x, y]w[x, y] = 0$ . Now choosing  $g(X) \in P_0$  such that  $[x, y]w = g(w[x, y])$ , we get  $k \in Z$  such that  $[x, y]w = kw[x, y]$ . Thus  $[x, y] \in C_l(R)$ , so by Lemma 2,  $[x, y] \in C$ .

THEOREM 4. *Let  $R$  be a periodic ring — that is, a ring in which for each  $x \in R$ , there exist distinct positive integers  $n, m$  for which  $x^n = x^m$ . If  $R$  is quasi-commutative and nilpotent elements commute with each other, then  $R$  is commutative.*

PROOF. Note first that in any periodic ring in which nilpotent elements commute, the set  $N$  of nilpotent elements is an ideal [2, Theorem 1]. Observe also that in any quasi-commutative  $R$ , idempotents are central. To see this, let  $e^2 = e$  and for a given  $x \in R$ , let  $f(X) \in Z[X]$  be such that  $e(ex - exe) = f((ex - exe)e)$ . Since the right side of this equation is either 0 or  $k \cdot 1$ , depending on whether  $R$  has 1, we have  $e(ex - exe)$  commuting with  $e$ . Thus,  $ex - exe = 0$ ; and similarly,  $xe - exe = 0$ .

Now consider  $R$  satisfying the hypotheses of Theorem 4, and express it as a subdirect product of subdirectly irreducible rings  $R_\alpha$ . Fixing  $\alpha$  and noting that in periodic rings, nilpotent elements of homomorphic images may be lifted to  $R$  [2, Proposition 0(c)], we see that  $R_\alpha$  inherits our original hypotheses. Since each element of  $R_\alpha$  has an idempotent power [2, Proposition 0(a)], either  $R_\alpha$  is nil, hence commutative, or  $R_\alpha$  has a non-zero idempotent. In the latter case, the subdirect irreducibility and the centrality of idempotents imply that  $R_\alpha$  has 1, hence by Lemma 3(a) satisfies the identity  $[[x, y], y] = 0$ . It follows, by a trivial modification of the third and fourth paragraphs of the proof of Theorem 2 of [3], that  $N \subseteq C$ . In a periodic ring, this condition guarantees commutativity [3, 7].

### 3. On commutativity of quasi-commutative rings with 1

The following theorem shows that any quasi-commutative ring with 1 is not too badly non-commutative.

THEOREM 5. *Let  $R$  be any quasi-commutative ring with 1. Then  $R$  satisfies the identity  $[[x, y], [z, w]] = 0$ ; moreover, if  $R$  is 3-torsion-free,  $R$  satisfies the identity  $[[x, y], w] = 0$ .*

PROOF. By Lemma 3(a),  $R$  satisfies the identity  $[[x, y], x] = 0$ , which can be linearized to yield the identity

$$(2) \quad [[x, y], w] + [x, [y, w]] = 0.$$

Substituting (2) several times into the Jacobi identity yields the identity

$$(3) \quad 3[[x, y], w] = 0,$$

which establishes the conclusion of the theorem when  $R$  is 3-torsion-free.

As usual, we need only consider the case of subdirectly irreducible  $R$ , which we may assume to possess 3-torsion and hence to be 2-torsion-free. Thus, if  $x \in R$  and  $x^2 = 0$ , the fact that  $[[x, y], x] = 0$  implies that  $xyx = 0$  for all  $y \in R$ . Since commutators square to zero by Lemma 3(b), applying the definition of quasi-commutativity

to  $[x, y]$  and  $w$  shows that for each  $x, y, w \in R$ , there exist  $k, m \in \mathbb{Z}$  such that

$$(4) \quad [x, y]w = k \cdot I + mw[x, y].$$

If  $R^+$  is not periodic — i.e., if  $I$  has infinite order in  $R^+$  — the fact that  $\mathcal{C}(R)$  is nil forces the  $k$  in (4) to be 0; thus commutators are in  $C_1(R)$ , hence in  $C$  by Lemma 2.

To complete our proof, therefore, we may suppose that there exists a positive integer  $q$  such that  $3^q R = 0$ . Moreover, we may assume  $q \geq 2$ , since the above argument may be modified to show that commutators are central if  $q = 1$ . In view of (2) and (3) and the fact that commutators square to zero, we need only show that if  $x$  is any element with  $x^2 = 0$  and  $3x \in C$ , then  $[x, y] \in C$  for all  $y \in R$ .

Consider such an element  $x$ , let  $y \in R$ , and let  $k_1, k_2, n_1, n_2 \in \mathbb{Z}$  be such that  $xy = k_1 \cdot I + n_1 yx$  and  $x(y+I) = k_2 \cdot I + n_2(y+I)x$ . Letting  $n_i - I = m_i$ ,  $i = 1, 2$ , and noting that  $[x, y] = [x, y+I]$ , we then have

$$(5) \quad [x, y] = k_1 \cdot I + m_1 yx$$

and

$$(6) \quad [x, y] = k_2 \cdot I + m_2(y+I)x.$$

If neither  $m_1$  nor  $m_2$  is divisible by 3, then (5) and (6), combined with the fact that  $[x, y], y] = 0$ , imply that both  $yx$  and  $(y+I)x$  commute with  $y$ , hence  $[x, y] = 0$ .

Let  $m_1 = 3^r M_1$  and  $m_2 = 3^s M_2$ , where  $(M_1, 3) = (M_2, 3) = 1$  and at least one of  $r, s$  is positive. If  $r \neq s$ , or if  $r = s$  and  $M_2 - M_1 \not\equiv 0 \pmod{3}$ , then subtracting (6) from (5) enables us to write  $3^r yx$  as an integral linear combination of  $I$  and  $3x$ , hence  $[x, y] \in C$  by (5). Thus, we may suppose that  $r = s$  and  $M_2 - M_1 \equiv 0 \pmod{3}$ . Rewriting (5) and (6), we have  $t, K_1, K_2 \in \mathbb{Z}$ ,  $t \geq 0$ , such that

$$(7) \quad [x, y] = 3^t K_1 \cdot I + 3^t M_1 yx$$

and

$$(8) \quad [x, y] = 3^t K_2 \cdot I + 3^t M_2(y+I)x,$$

where at least one of  $K_1$  and  $K_2$  is not divisible by 3. Subtracting (7) from (8) yields

$$(9) \quad 3^t (K_2 - K_1) \cdot I + 3^t (M_2 - M_1) yx + 3^t M_2 x = 0.$$

Now if  $3^t x = 0$ , it follows from (7) that  $[x, y] \in C$ , so we may assume that the order of  $x$  in  $R^+$  is  $3^{t+n}$  for some  $n \geq 1$ . Consider first the possibility that  $n > 1$ . Multiplying (7) by  $3^{t+n-1}$  and using the fact that  $3x \in C$ , we get

$$(10) \quad 3^{t+n-1} K_1 \cdot I + 3^{t+n-1} M_1 yx = 0;$$

multiplying (9) by  $3^{t+n-1}$  and recalling that  $M_2 - M_1 \equiv 0 \pmod{3}$ , we get

$$(11) \quad 3^{t+n-1} (K_2 - K_1) \cdot I + 3^{t+n-1} M_2 x = 0.$$

Assuming temporarily that  $(K_1, 3) = 1$ , we appeal to (10) to obtain  $3^{t+n} \cdot I = 0$ ; and it follows from (11) that  $K_2 - K_1 \not\equiv 0 \pmod{3}$ . Thus, comparing (10) and (11) shows that  $3^{t+n-1} x = \pm 3^{t+n-1} yx$ . On the other hand, if  $(K_1, 3) \neq 1$  but  $(K_2, 3) = 1$ , then we can simply repeat our entire argument with  $y+I$  and  $(y+I) - I$  in the roles of  $y$  and  $y+I$ , and obtain  $3^{t+n-1} x = \pm 3^{t+n-1} (y+I)x$ . Thus, for every  $i \in \mathbb{Z}$ , an



appropriate choice of  $w$  as  $y+i$  or  $y+i+1$  yields

$$(12) \quad 3^{r+n-1}x = \pm 3^{r+n-1}wx.$$

It follows that there are distinct  $i_1, i_2 \in Z$  with  $|i_1 - i_2| = 2$  such that  $w_1 = y + i_1$  and  $w_2 = y + i_2$  both satisfy (12); and adding or subtracting the two versions of (12) yields the contradiction  $3^{r+n-1}x = 0$ .

What remains to be considered is the case  $n=1$ . Assuming without loss that  $(K_1, 3)=1$  in equation (7), we see from (7) that  $3^{t+1} \cdot 1 = 0$ , so that  $t \geq q-1$ ; and from (9) and the previous assumption that  $M_2 - M_1 \equiv 0 \pmod{3}$ , it follows that  $3^r x$ , and hence  $3^r yx$ , belongs to  $3^{q-1}R$ . Thus, in view of (7), we are finished once we show that  $3^{q-1}R \subseteq C$ .

Accordingly, let  $x$  and  $y$  be arbitrary elements in  $R$ . By repeating earlier arguments, we may assume there are  $j, k \in Z$  such that  $[3^{q-1}x, y] = j \cdot 1 + ky(3^{q-1}x)$  and  $k \equiv 0 \pmod{3}$ ; therefore,  $[3^{q-1}x, y] = j \cdot 1$ . Multiplying this equality by 3 shows that  $3j \cdot 1 = 0$ , so that  $3^{q-1} | j$ . If  $3^q | j$ , we have  $[3^{q-1}x, y] = 0$ , which is what we want. Otherwise  $3^{q-1}(1 \pm [x, y]) = 0$ , which implies the contradiction  $3^{q-1} \cdot 1 = 0$ , since  $1 + [x, y]$  and  $1 - [x, y]$  are each invertible. This completes the proof of Theorem 5.

#### 4. Certain $P$ -commutative rings are commutative

Chacron and Thierrin [5] have shown that if  $P$  consists of  $X$  together with all polynomials of co-degree at least two, then any  $P$ -commutative ring  $R$  must be commutative. In general, no significant improvement on this result can be expected; however, as the following theorems show, if  $R$  has  $1$ , the situation is much better.

**THEOREM 6.** *Let  $P_1$  be the set of all polynomials in  $Z[X]$  which are either of co-degree at least two or of form  $nX$  for some  $n \in Z$ . Then every  $P_1$ -commutative ring  $R$  with  $1$  is commutative.*

**PROOF.** Recall from the proof of Theorem 4 that all idempotents are central. Moreover, assume henceforth that  $R$  is subdirectly irreducible, in which case  $1$  is the only non-zero idempotent.

Suppose now that  $x \in R$  has right inverse  $y$ . Then  $xy = 1$  and  $yx = k \cdot 1$  for some  $k \in Z$ , so that

$$(13) \quad [x, y] = (1 - k) \cdot 1.$$

Since  $[x, y]$  is nilpotent by Theorem 1, we see at once that  $k=1$  in the case where  $1$  has infinite order in  $R^+$ . On the other hand, if  $1$  has finite order in  $R^+$ , the subdirect irreducibility guarantees a prime  $p$  and a positive integer  $r$  for which  $p^r \cdot 1 = 0$ ; and it follows from (13) and the nilpotence of  $[x, y]$  that  $k \equiv 1 \pmod{p}$ , so that for an appropriate exponent  $n$  the Euler—Fermat Theorem gives  $(yx)^n = 1$ . Hence any element having a one-sided inverse is invertible. An obvious remark of the same kind is that  $xy=0$  if and only if  $yx=0$ .

Suppose now that  $R$  is a (subdirectly irreducible) counterexample, and that  $y \notin C$ . By Lemma 2, there must exist  $u \in R$  such that  $yu \neq kuy$  for all  $k \in Z$ . It follows that there exists a polynomial  $f(X)$  of co-degree at least two such that  $yu = f(uy)$ ; and since  $uy = g(yu)$  for some polynomial  $g(X)$  of co-degree at least one,

we get a polynomial  $h(X) = X - X^2 h_1(X)$  for which  $h(yu) = 0$ . Since  $yu = (yu)^2 h_1(yu)$ , it follows that  $yuh_1(yu)$  is a non-zero idempotent, necessarily 1; hence  $u$  is invertible. It is immediate that  $u^{-1}yu \neq ku(u^{-1}y)$  for each  $k \in Z$ , so the above argument gives a polynomial  $f_0(X) = X - X^2 f_1(X)$  for which  $0 = f_0(u^{-1}yu) = u^{-1}f_0(y)u = f_0(y)$ . Thus, for every  $x \in R$ , there exists a polynomial  $g_x(X) \in Z[X]$  such that  $x - x^2 g_x(x) \in C$ .

At this point we could conclude that  $R$  is commutative by appealing to a well-known theorem of Herstein [6]; however, in order to produce a more elementary proof and to allow for subsequent generalization, we take a different approach. Note that by the argument above, we know that all non-invertible elements — in particular, all zero divisors — are central. Moreover, since all non-trivial one-sided annihilators are two-sided, it is easy to show that the set  $D$  of zero divisors is the annihilator of the heart of  $R$ ; hence,  $D$  is an ideal.

Now suppose  $x \notin C$  and  $[x, y] \neq 0$ . Let  $k(X) \in Z[X]$  be a polynomial of smallest positive degree for which  $k(x) \in C$ . Since  $[[x, y], x] = 0$ , the statement that  $[k(x), y] = 0$  may be written as  $k'(x)[x, y] = 0$ , where  $k'(X)$  is the formal derivative of  $k(X)$ . Because  $k'(x) \in D \subseteq C$ ,  $k'(X)$  must be a non-zero constant polynomial; thus,  $\bar{R} = R/D$  has non-zero characteristic  $p$ . Moreover, if  $(n, p) = 1$  and  $k_1(X)$  is a polynomial of form  $X^n - X^{2n}g(X^n)$  for which  $k_1(x) \in C$ , the result that  $k'_1(x)[x, y] = 0$  shows that in  $\bar{R}$  the element  $\bar{x} = x + D$  is algebraic over the prime subfield  $\bar{P}$ ; hence  $\bar{x}$  generates a finite subfield  $\bar{F}$ . If  $\bar{F}$  has  $p^s$  elements, then

$$(14) \quad x - x^{p^s} \in D \subseteq C.$$

On the other hand, since  $\bar{R}$  has characteristic  $p$ , we have  $px \in D \subseteq C$ ; hence  $p[x, y] = 0$  for all  $y \in R$ , and therefore  $[x^p, y] = px^{p-1}[x, y] = 0$  for all  $y \in R$ . Thus  $x^p \in C$ , which combines with (14) to contradict our assumption that  $x \notin C$ . Consequently,  $R$  must be commutative.

In fact, we can manage with a larger set than  $P_1$ , as the next theorem shows.

**THEOREM 7.** *Let  $P_2$  be the subset of  $Z[X]$  containing all constant polynomials, all polynomials of form  $nX$ , and all polynomials of co-degree at least two. Then every  $P_2$ -commutative ring with 1 is commutative.*

**PROOF.** Let  $P_1$  be as in Theorem 6, and define the unordered pair  $\{x, y\}$  of elements of  $R$  to be a  $P_1$ -pair if there exist  $f(X)$  and  $g(X)$  in  $P_1$  for which  $xy = f(yx)$  and  $yx = g(xy)$ . Note that if  $R$  has 1 and all pairs  $\{x, y\}$  are  $P_1$ -pairs, then  $R$  is commutative by Theorem 6.

Suppose that  $R$  is a counterexample to Theorem 7, and that  $\{x, y\}$  is not a  $P_1$ -pair. Then  $xy = j \cdot 1$  and  $yx = k \cdot 1$  for appropriate  $j, k \in Z$ , hence  $[x, y] = (j - k) \cdot 1$ . If 1 had infinite order in  $R^+$  or if there were a prime  $p$  for which  $pR = 0$ , then the nilpotence of  $[x, y]$  would force  $j \cdot 1 = k \cdot 1$ , contrary to the supposition that  $[x, y] \neq 0$ . Thus, by replacing  $R$  by a direct summand if necessary, we may assume that there exist a prime  $p$  and an integer  $r > 1$  for which  $p^r R = 0$ . We may further assume that  $R$  has been chosen so that  $r$  is minimal among all such counterexamples, and that  $R$  is subdirectly irreducible.

Letting  $I = \{x \in R \mid px = 0\}$ , note that the factor ring  $\bar{R} = R/I$  has  $p^{r-1}\bar{R} = 0$ , hence is commutative; thus,  $p[x, y] = 0$  for all  $x, y \in R$ . Note also that if  $y$  is an element of  $R$  for which  $\{y, x\}$  is a  $P_1$ -pair for all  $x \in R$ , then a piece of the proof of Theorem 6 shows that  $y$  must be either central or invertible.

Choose an unordered pair  $\{y, u\}$  which is not a  $P_1$ -pair, such that  $y$  has minimal additive order, say  $p^s$ , among all elements belonging to a non- $P_1$ -pair. Assume without loss that there exists no  $f(X) \in P_1$  for which  $yu = f(uy)$ , so that there exist,  $j, k \in \mathbb{Z}$  such that  $uy = j \cdot I$  and  $yu = k \cdot I \neq 0$ . Since  $(k-j)y = (yu)y - y(uy) = 0$ , we have  $k \equiv j \pmod{p^s}$ , from which we see that  $r = s$  would imply the contradiction  $[y, u] = 0$ . Therefore  $s < r$ ; in particular,  $y$  cannot be invertible.

Suppose first that there exists a non-negative integer  $t < r$  for which  $yu = ap^t \cdot I$  and  $uy = bp^t \cdot I$ , where  $(a, p) = (b, p) = 1$ . The congruence  $nb \equiv a \pmod{p^{r-t}}$  has a solution  $n_0$ , and it follows that  $yu = n_0 uy$ , contradicting our choice of  $y$ .

Next consider the possibility that there exist  $a, b \in \mathbb{Z}$  with  $(a, p) = (b, p) = 1$ , and distinct non-negative integers  $t$  and  $T$ , both less than  $r$ , for which  $yu = ap^t \cdot I$  and  $uy = bp^T \cdot I$ . Then  $0 = p[y, u] = ap^{t+1} \cdot I - bp^{T+1} \cdot I$ , and multiplying by the smaller of  $p^{r-t-1}$  and  $p^{r-T-1}$  yields the contradiction that  $p^{r-1} \cdot I = 0$ .

In view of the facts that  $p[y, u] = 0$  and  $yu \neq 0$ , the only remaining possibility is that  $yu = ap^{r-1} \cdot I$  and  $uy = 0$ , where  $(a, p) = 1$ . Since  $[y, u] \neq 0$ , neither  $(u+I)y$  nor  $y(u+I)$  can be of form  $t \cdot I$ , hence  $\{y, u+I\}$  is a  $P_1$ -pair. Moreover, since a straightforward computation yields  $(y(u+I))^k = y^k$  for all  $k \geq 2$ , the existence of  $f(X)$  of co-degree at least 2 for which  $(u+I)y = f(y(u+I))$  would imply that  $y = f(y)$ , and hence that  $y$  is invertible; therefore, we must have  $n \in \mathbb{Z}$  such that  $(u+I)y = ny(u+I)$ , from which it follows that

$$(15) \quad (n-1)y = -nyu = -nap^{r-1} \cdot I.$$

Now if  $n \not\equiv 1 \pmod{p}$ , we can express  $y$  as an integral multiple of  $I$ , contrary to the fact that  $[y, u] \neq 0$ ; therefore, let  $n = 1 + bp^t$  with  $t \geq 1$  and  $(b, p) = 1$ . Recall that  $y$  has additive order  $p^s$ . If  $t \geq s$ , which is certainly the case if  $s = 1$ , (15) implies the impossible result  $p^{r-1} \cdot I = 0$ ; on the other hand, if  $t \leq s-2$ , multiplying (15) by an appropriate power of  $p$  yields  $p^{s-1}y = 0$ , again an impossibility. We are left with  $t = s-1$ , and (15) gives  $bp^{s-1}y = -nap^{r-1} \cdot I$ . But this implies that  $w = by + nap^{r-s} \cdot I$  has additive order less than  $p^s$ , so that every pair  $\{w, x\}$  is a  $P_1$ -pair; and since  $w$  cannot be invertible,  $w \in C$  and hence  $y \in C$ . This final contradiction demolishes the assumption that  $R$  was a counterexample, thereby completing the proof of Theorem 7.

**COROLLARY 8.** *Let  $R$  be a ring with 1; and suppose that for each  $x, y \in R$ , there exist  $k, n \in \mathbb{Z}$  with  $n \geq 0$ , for which  $xy = k(yx)^n$ . Then  $R$  is commutative.*

## 5. Some extensions

It is natural to consider what happens if  $R$  is an algebra over a field  $F$  and the polynomials used in the definition of quasi-commutativity and the quasi-centers are in  $F[X]$ . In fact, Lemma 2 and Theorems 1—4 remain true for arbitrary  $F$ ; and Theorems 6 and 7 certainly hold if  $F$  is either finite or of characteristic 0. Straightforward modifications — usually simplifications — of the given proofs make these assertions clear.

## REFERENCES

- [1] BELL, H. E., On some commutativity theorems of Herstein, *Arch. Math. (Basel)*, **24** (1973), 34—38. *MR* **47** # 8631.
- [2] BELL, H. E., Some commutativity results for periodic rings, *Acta Math. Acad. Sci. Hungar.*, **28** (1976), 279—283. *MR* **54** # 7556.
- [3] BELL, H. E., On commutativity of periodic rings and near-rings, *Acta Math. Acad. Sci. Hungar.* **36** (1980), 293—302. *MR* **82h**: 16025.
- [4] BELL, H. E., Quasi-centers, quasi-commutators, and ring commutativity, *Acta Math. Hungar.* **41** (1983), 127—136.
- [5] CHACRON, M. and THIERRIN, G.,  $\sigma$ -reflexive semigroups and rings, *Canad. Math. Bull.*, **15** (1972), 185—188. *MR* **46** # 5487.
- [6] HERSTEIN, I. N., The structure of a certain class of rings, *Amer. J. Math.*, **75** (1953), 864—871. *MR* **15**—392.
- [7] HERSTEIN, I. N., A note on rings with central nilpotent elements, *Proc. Amer. Math. Soc.*, **5** (1954), 620. *MR* **16**—5.

(Received July 20, 1983)

DEPARTMENT OF MATHEMATICS  
BROCK UNIVERSITY  
ST. CATHARINES, ONTARIO  
L2S 3A1  
CANADA

# AN INEQUALITY IN COMPLEX RIESZ ALGEBRAS

C. B. HUIJSMANS

A (real) Riesz algebra (lattice ordered algebra)  $A$  is, by definition, an associative algebra which is at the same time a Riesz space (vector lattice) with the additional property that the multiplication and the ordering are compatible, that is

$$u, v \in A^+ \Rightarrow uv \in A^+,$$

where

$$A^+ = \{f \in A: f \geq 0\}.$$

We mention two immediate consequences of the definition.

1)  $u, v \in A, w \in A^+$  and  $u \geq v$  implies  $uw \geq vw$ .

2)  $|fg| \leq |f||g|$  for all  $f, g \in A$ .

For the second property, it follows from  $fg = f^+g^+ + f^-g^- - (f^+g^- + f^-g^+)$  that  $fg \leq f^+g^+ + f^-g^-$ . Since  $f^+g^+ + f^-g^- \geq 0$  as well, we find

$$(fg)^+ \leq f^+g^+ + f^-g^-$$

and similarly

$$(fg)^- \leq f^+g^- + f^-g^+.$$

Hence,

$$|fg| = (fg)^+ + (fg)^- \leq (f^+ + f^-)(g^+ + g^-) = |f||g|.$$

The question arises immediately whether the corresponding inequality also holds in complex Riesz algebras.

We recall some facts. First of all, we make from now on the blanket assumption that all Riesz spaces and all Riesz algebras we consider are Archimedean and relatively uniformly complete (for elementary Riesz space theory and terminology we refer to [3], [5] and [8]).

Let  $L$  be an Archimedean relatively uniformly complete Riesz space. As is well-known, one can introduce in the vector space complexification  $L + iL$  of  $L$  a modulus defined by ( $\varphi = f + ig$ ,  $f = \operatorname{Re} \varphi$ ,  $g = \operatorname{Im} \varphi$ ):

$$\begin{aligned} |\varphi| &= |f + ig| = \sup \{ \operatorname{Re} (e^{-i\theta} \varphi), 0 \leq \theta \leq 2\pi \} \\ &= \sup \{ f \cos \theta + g \sin \theta, 0 \leq \theta \leq 2\pi \} \\ &= \sup \{ |f \cos \theta + g \sin \theta|, 0 \leq \theta \leq 2\pi \} \end{aligned}$$



(see [4], section 3). The modulus satisfies the following properties:

- (a)  $|\varphi + \psi| \leq |\varphi| + |\psi|$  for all  $\varphi, \psi \in L + iL$ .
- (b)  $||\varphi| - |\psi|| \leq |\varphi + \psi|$  for all  $\varphi, \psi \in L + iL$ .
- (c)  $|\alpha\varphi| = |\alpha| |\varphi|$  for all  $\alpha \in \mathbb{C}$ ,  $\varphi \in L + iL$ .
- (d)  $|\operatorname{Re} \varphi| \leq |\varphi|$ ,  $|\operatorname{Im} \varphi| \leq |\varphi|$  for all  $\varphi \in L + iL$ .

It is common to call  $L + iL$ , equipped with the above modulus, a complex Riesz space.

Now, let  $A$  be an Archimedean relatively uniformly complete Riesz algebra. The algebra complexification  $A + iA$  is a complex Riesz space as well. Of course, it is very natural to ask whether the multiplicative analogue of the complex triangle inequality (a)

$$|\varphi\psi| \leq |\varphi| |\psi|$$

holds for all  $\varphi, \psi \in A + iA$ . The main purpose of the present note is to show the validity of this inequality, that, somewhat surprisingly, does not seem to occur in the literature. However, in some special cases the inequality is known to be true. We mention them.

First, if  $A$  is a semiprime  $f$ -algebra, then even equality holds, that is

$$|\varphi\psi| = |\varphi| |\psi|$$

holds for all  $\varphi, \psi \in A + iA$ . The proof is mainly based on the fact that now  $|\varphi|^2 = f^2 + g^2$  if  $\varphi = f + ig$  ([2], section 5).

There is another important class of complex Riesz algebras in which the above inequality is known to be true, the so-called complex Banach lattice algebras. Recall that a real Banach lattice  $A$  which is simultaneously a Riesz algebra is said to be a real Banach lattice algebra whenever  $\|uv\| \leq \|u\| \|v\|$  for all  $u, v \in A^+$ . It is not difficult to prove that  $A$  is a real Banach algebra. Since any real Banach lattice is both Archimedean and relatively uniformly complete, the algebra complexification  $A + iA$  of  $A$  is a complex Riesz space. Equipped with the norm  $\|\varphi\| = \| |\varphi| \|$  (by definition)  $A + iA$  is a so-called complex Banach lattice. It turns out that  $A + iA$  is a complex Banach algebra in which  $|\varphi\psi| \leq |\varphi| |\psi|$  holds for all  $\varphi, \psi \in A + iA$ . For details and proofs we refer the reader to the thesis of W. Arendt ([1], Lemma 1.5) and the paper by E. Scheffold ([6], Theorem 1.1).

Finally, we mention the Riesz algebra  $\mathcal{L}_b(L)$  of all order bounded operators on some Dedekind complete Riesz space  $L$ . The modulus

$$|T| = \sup \{T_1 \cos \theta + T_2 \sin \theta, 0 \leq \theta \leq 2\pi\}$$

of an order bounded operator  $T = T_1 + iT_2$  on the complexification  $L + iL$  of  $L$  satisfies

$$|T|u = \sup \{ |T\varphi| : \varphi \in L + iL, |\varphi| \leq u \}$$

for all  $u \in L^+$ . Consequently,  $|T\varphi| \leq |T| |\varphi|$  for all  $\varphi \in L + iL$ . These results are, independently, due to both H. H. Schaefer ([5], chapter 4, Theorem 1.8) and W. J. de Schipper ([7], Corollary 3.5). The complexification of the Dedekind complete Riesz algebra  $\mathcal{L}_b(L)$  is  $\mathcal{L}_b(L + iL)$ . It follows immediately from the above formulas

that

$$|ST| \leq |S||T|$$

for all  $S, T \in \mathcal{L}_b(L+iL)$ . Indeed, if  $u \in L^+$ , then we have

$$|ST\varphi| \leq |S||T\varphi| \leq |S||T||\varphi| \leq |S||T|u$$

for all  $\varphi \in L+iL$  with  $|\varphi| \leq u$ . Hence,

$$|ST|u = \sup \{ |ST\varphi| : \varphi \in L+iL, |\varphi| \leq u \} \leq |S||T|u.$$

This holds for all  $u \in L^+$ , whence  $|ST| \leq |S||T|$ .

So far the known results. For the proof of the general case we need an approximation lemma which can be proved by means of Yosida's representation theorem and the use of a suitable partition of unity on a compact Hausdorff space. For details we refer to [5], chapter 4, Theorem 1.8 and [7], Proposition 2.7.

**LEMMA 1** (the approximation property). *Let  $L$  be an Archimedean relatively uniformly complete Riesz space, let  $u \in L^+$ ,  $\varphi \in L+iL$  with  $|\varphi| \leq u$  and  $\varepsilon > 0$ . Then there exist  $\alpha_1, \dots, \alpha_n \in \mathbb{C}$ ,  $v_1, \dots, v_n \in L^+$  such that*

$$(I) \quad \sum_{k=1}^n v_k = u,$$

$$(II) \quad |\alpha_k| \leq 1 \quad (k=1, \dots, n),$$

$$(III) \quad \left| \varphi - \sum_{k=1}^n \alpha_k v_k \right| \leq \varepsilon u.$$

**THEOREM 2.** *Let  $A$  be an Archimedean relatively uniformly complete Riesz algebra. Then*

$$|\varphi\psi| \leq |\varphi||\psi|$$

for all  $\varphi, \psi \in A+iA$ .

**PROOF.** Put  $\varphi = f+ig$ ,  $\psi = h+ik$ . For ease of survey we divide the proof in several steps.

*Step 1.*  $\psi \in A$  (i.e.,  $k=0$ ). In this case  $\varphi\psi = fh+i(gh)$ , so

$$\begin{aligned} |\varphi\psi| &= \sup \{ fh \cos \theta + gh \sin \theta, 0 \leq \theta \leq 2\pi \} \\ &= \sup \{ |fh \cos \theta + gh \sin \theta|, 0 \leq \theta \leq 2\pi \}. \end{aligned}$$

Further,

$$\begin{aligned} |fh \cos \theta + gh \sin \theta| &\leq |f \cos \theta + g \sin \theta| |h| \leq \\ &\leq \sup \{ |f \cos \theta + g \sin \theta|, 0 \leq \theta \leq 2\pi \} |h| = |\varphi| |h| = |\varphi| |\psi| \end{aligned}$$

for all  $\theta$ , where we use properties 1) and 2) of the real Riesz algebra. By taking the supremum on the left-hand side we find  $|\varphi\psi| \leq |\varphi||\psi|$ .

*Step 2.*  $|\varphi\psi| \leq 2|\varphi||\psi|$  for all  $\varphi, \psi \in A+iA$ . Indeed,

$$|\varphi\psi| = |\varphi h + i\varphi k| \leq |\varphi h| + |\varphi k| \leq |\varphi| |h| + |\varphi| |k| \leq 2|\varphi| |\psi|,$$

where we use step 1 and properties (a), (c) and (d) of the modulus.

*Step 3.* By the approximation property there exists a sequence  $\{\varphi_n\}_{n=1}^\infty$  in  $A+iA$  such that

$$|\varphi - \varphi_n| \leq n^{-1}|\varphi| \quad (n = 1, 2, \dots),$$

where each  $\varphi_n$  is a finite sum of the form  $\sum \alpha_k v_k$  with  $\sum |\alpha_k| v_k \leq |\varphi|$ . Observe now that

$$||\varphi\psi| - |\varphi_n\psi|| \leq |(\varphi - \varphi_n)\psi| \leq 2|\varphi - \varphi_n||\psi| \leq 2n^{-1}|\varphi||\psi| \quad (n = 1, 2, \dots),$$

where we use step 2, property 1) of the Riesz algebras and property (b) of the modulus. Therefore, once we can show that  $|\varphi_n\psi| \leq |\varphi||\psi|$  ( $n=1, 2, \dots$ ), the inequality

$$||\varphi\psi| - |\varphi_n\psi|| \leq 2n^{-1}|\varphi||\psi| \quad (n = 1, 2, \dots)$$

implies directly that  $|\varphi\psi| \leq |\varphi||\psi|$ . The proof is straightforward. Indeed,

$$|\varphi_n\psi| = |\sum \alpha_k v_k \psi| \leq \sum |\alpha_k| |v_k \psi| \leq \sum |\alpha_k| v_k |\psi| \leq |\varphi||\psi|,$$

where we use property (a) and (c) of the modulus, step 1 and property 1) of the real Riesz algebras. The proof is complete.

#### REFERENCES

- [1] ARENDT, W., Über das Spektrum regulärer Operatoren, Dissertation, Tübingen, 1979.
- [2] BEUKERS, F., HUIJSMANS, C. B. and PAGTER, B. DE, Unital embedding and complexification of  $f$ -algebras, *Math. Z.* **183** (1983), 131—144. *MR* **85c**: 06016.
- [3] LUXEMBURG, W. A. J. and ZAAENEN, A. C., *Riesz spaces, Vol. I*, North-Holland Mathematical Library, North-Holland Publishing Co., Amsterdam—London—New York, 1971. *MR* **58** # 23483.
- [4] LUXEMBURG, W. A. J. and ZAAENEN, A. C., The linear modulus of an order bounded linear transformation I, *Nederl. Akad. Wetensch. Proc. Ser. A.=Indag. Math.* **33** (1971), 422—434. *MR* **46** # 2475a.
- [5] SCHAEFER, H. H., *Banach lattices and positive operators*, Grundlehren der mathematischen Wissenschaften Band **215**, Springer-Verlag, New York—Heidelberg, 1974. *MR* **54** # 11023.
- [6] SCHEFOLD, E., Über komplexe Banachverbandsalgebren, *J. Funct. Anal.* **37** (1980), 382—400. *MR* **81m**: 46064.
- [7] DE SCHIPPER, W. J., A note on the modulus of an order bounded linear operator between complex vector lattices, *Nederl. Akad. Wetensch. Proc. Ser. A* **76** = *Indag. Math.* **35** (1973), 355—367. *MR* **51** # 13757.
- [8] ZAAENEN, A. C., *Riesz spaces, Vol. II*, North-Holland Mathematical Library, North-Holland Publishing Co., Amsterdam—New York—Oxford, 1983. *MR* **86b**: 46001.

(Received July 26, 1982)

MATHEMATISCH INSTITUUT  
WASSENAARSEWEG 80  
NL—2333 AL LEIDEN  
THE NETHERLANDS

# BOUNDING THE SIZE OF A FAMILY KNOWING THE CARDINALITY OF DIFFERENCES

P. FRANKL

## Abstract

Suppose  $\mathcal{F}$  is a family of distinct subsets of an  $n$ -element set and for  $F, F'$  in the family  $|F - F'|$  takes up at most  $s$  non-zero values. Then  $|\mathcal{F}| \leq \sum_{0 \leq i \leq s} \binom{n}{i}$ . Equality holds if one takes all subsets of size at most  $s$ . We prove this in a stronger, mod  $p$  form, and extend a result of Delsarte, too.

## 1. Introduction

Let  $X$  be an  $n$ -element set and  $\mathcal{F}$  a family of its subsets, i.e.  $\mathcal{F} \subseteq 2^X$ . Katona [4] asked to find upper bounds for  $|\mathcal{F}|$  assuming  $|F - F'| \leq s$  holds for all  $F, F' \in \mathcal{F}$ . Here we consider the following more general situation. Suppose we are given a prime  $p$  and  $s$  distinct non-zero residues modulo  $p: r_1, \dots, r_s$  so that

$$(1) \quad |F - F'| \equiv r_1 \text{ or } \dots \text{ or } r_s \pmod{p} \quad \text{for all } F \not\subseteq F', F, F' \in \mathcal{F}.$$

THEOREM 1. *For a family  $\mathcal{F}$  satisfying (1) we have*

$$(2) \quad |\mathcal{F}| \leq \sum_{0 \leq i \leq s} \binom{n}{i}.$$

Note that in the special case when all sets in  $\mathcal{F}$  have the same size  $k$  one has  $|F \cap F'| = |F - F'|$ . Consequently, in this case, (1) is equivalent to: all pairwise intersections lie in  $s$  distinct residue classes, different from the class of  $k$ , modulo  $p$ .

In [3] it was proved that this stronger condition implies  $|\mathcal{F}| \leq \binom{n}{s}$ .

Choosing  $p > n$  we obtain:

COROLLARY 1. *If  $|F - F'|$  takes only  $s$  distinct non-zero values for  $F, F' \in \mathcal{F}$ , then (2) holds.*

Specifying to  $r_i = i$ ,  $1 \leq i \leq s$ , we deduce:

COROLLARY 2. *If  $|F - F'| \leq s$  holds for all  $F, F' \in \mathcal{F}$  then (2) holds.*

This corollary solves the problem of Katona since one has equality in (2) by taking  $\mathcal{F} = \{F \subset X: |F| \leq s\}$ .

1980 *Mathematics Subject Classification*. Primary 05A05.  
*Key words and phrases*. Subset, matrix, rank.

## 2. The proof of Theorem 1

We may suppose without loss of generality  $s < n$ . For  $0 \leq i \leq s$  let us define:

$$\mathcal{F}^{(i)} = \{G: |G| = i, G \subseteq F \text{ holds for some } F \in \mathcal{F}\}, \quad f_i = |\mathcal{F}^{(i)}|, \quad |\mathcal{F}| = f.$$

Now we want to define the  $i$ 'th containment and disjointness matrices:  $C_i$  and  $D_i$ , resp. For that let us order the elements of  $\mathcal{F}$  and  $\mathcal{F}^{(i)}$  in an arbitrary way:

$$\mathcal{F}^{(i)} = \{G_1^i, \dots, G_{f_i}^i\}, \quad 0 \leq i \leq s, \quad \mathcal{F} = \{F_1, \dots, F_f\}, \quad |F_1| \leq |F_2| \leq \dots \leq |F_f|.$$

The matrices  $C_i, D_i$  are  $f_i$  by  $f$  matrices with respective general entry  $c_i(a, b)$ ,  $d_i(a, b)$  defined by:

$$c_i(a, b) = \begin{cases} 0 & \text{if } G_a^i \not\subseteq F_b \\ 1 & \text{if } G_a^i \subseteq F_b; \end{cases} \quad d_i(a, b) = \begin{cases} 0 & \text{if } G_a^i \cap F_b \neq \emptyset \\ 1 & \text{if } G_a^i \cap F_b = \emptyset. \end{cases}$$

Obviously  $\text{rank } (D_i) \leq f_i \leq \binom{n}{i}$  holds. Thus the row space, say  $V_i$ , over the rationals has dimension at most  $f_i$ .

Our aim is to find a non-singular  $f$  by  $f$  matrix,  $M$  whose row space is contained in the span of  $V_0, \dots, V_s$ . In fact, since this latter has dimension at most  $f_0 + f_1 + \dots + f_s$  while the row space of  $M$  has full dimension  $f$ , the inequality  $f \leq f_0 + \dots + f_s$  and thus the statement of Theorem 1 follows.

Denoting by  $C_i^*$  the transposed of  $C_i$  let us calculate the general entry  $i(a, b)$  of the  $f$  by  $f$  matrix  $C_i^* D_i$ .

Since both  $C_i$  and  $D_i$  are  $(0, 1)$ -matrices  $i(a, b)$  is simply the number of  $G \in \mathcal{F}^{(i)}$  satisfying  $G \subseteq F_a$ ,  $G \cap F_b = \emptyset$  i.e.,  $G \subseteq F_a - F_b$ . Hence

$$(3) \quad i(a, b) = \binom{|F_a - F_b|}{i} \quad \text{holds.}$$

Note that the row space of  $C_i^* D_i$  is contained in that of  $D_i$ . Let us define  $q(x) = (x - r_1)(x - r_2) \dots (x - r_s)$ . Being a polynomial of degree  $s$  with integer coefficients,  $q(x)$  can be uniquely expressed in the form

$$(4) \quad q(x) = \sum_{0 \leq i \leq s} a_i \binom{x}{i}$$

(where  $a_0, \dots, a_s$  are integers).

Let us define  $M = \sum_{0 \leq i \leq s} a_i C_i^* D_i$ .

Now the row space of  $M$  is contained in the span of  $V_0, V_1, \dots, V_s$  (the row spaces of  $D_0, \dots, D_s$ ). Thus it is sufficient to show that  $M$  has full rank:  $\text{rank } M = f$ .

In view of (3) and (4) the matrix  $M$  has general entry

$$m(a, b) = (|F_a - F_b| - r_1) \dots (|F_a - F_b| - r_s), \quad 1 \leq a, b \leq f.$$

If  $a > b$ , then  $F_a \not\subseteq F_b$ , whence  $|F_a - F_b| \equiv r_1$  or  $\dots$  or  $r_s \pmod{p}$ . Consequently  $m(a, b)$  is zero modulo  $p$ . This means that all the entries below the diagonal of  $M$  are zero mod  $p$ . We infer

$$\det M \equiv \prod_{1 \leq a \leq f} m(a, a) = ((-1)^s r_1 r_2 \dots r_s)^f \not\equiv 0 \pmod{p}$$



### 3. Symmetric differences

For the case of symmetric difference the inequality corresponding to Corollary 1 was proved by Delsarte [1]. Actually his proof can be readapted to yield the following stronger version:

**THEOREM 2.** Suppose that  $r_1, \dots, r_s$  are distinct non-zero residues modulo a prime  $p$  so that for any pair of distinct  $F, F' \in \mathcal{F}$

$$(5) \quad |F \nabla F'| \equiv r_1 \text{ or } \dots \text{ or } r_s \pmod{p}.$$

Then (2) holds.

**PROOF.** Fix an arbitrary ordering of the  $i$ -subsets of  $X$ ,  $0 \leq i \leq s$ . Define the  $\binom{n}{i}$  by  $f$  matrix  $H_i$  by  $(G_i^t$  is the  $t$ 'th  $i$ -subset)

$$h_i(t, a) = (-1)^{|F_a \cap G_i^t|}.$$

Defining  $N_i = H_i^* H_i$  we have  $\text{rank}(H_i^* H_i) = \text{rank } H_i \leq \binom{n}{i}$ , and  $N_i$  has general entry

$$\begin{aligned} n_i(a, b) &= \sum_t (-1)^{|F_a \cap G_i^t|} (-1)^{|F_b \cap G_i^t|} = \sum_t (-1)^{|(F_a \nabla F_b) \cap G_i^t|} = \\ &= \sum_{0 \leq j \leq i} (-1)^j \binom{|F_a \nabla F_b|}{j} \binom{n - |F_a \nabla F_b|}{i - j}. \end{aligned}$$

Thus  $n_i(a, b)$  is a polynomial of degree  $i$  in  $|F_a \nabla F_b|$ . Let us denote this polynomial by  $p_i(x)$ , actually it is the Krawtchouk polynomial of degree  $i$  (cf. Delsarte [2]). Now the proof can be finished as in the case of Theorem 1. Define real numbers  $b_i$ ,  $0 \leq i \leq s$ , by  $\sum_{0 \leq i \leq s} b_i p_i(x) = (x - r_1) \dots (x - r_s)$  and consider the matrix  $N = \sum_{0 \leq i \leq s} b_i N_i$ .

All off-diagonal entries of  $N$  are zero modulo  $p$  but none of the diagonal entries.

Thus  $N$  has full rank:  $f$ , yielding  $f \leq \sum_{0 \leq i \leq s} \text{rank } N_i \leq \sum_{0 \leq i \leq s} \binom{n}{i}$   $\square$ .

Let us note that this theorem was proved independently by A. Blokhuis.

### 4. Open problems

**PROBLEM 1.** Suppose  $\mathcal{F}$  is an antichain (i.e. no member of  $\mathcal{F}$  contains an other one) satisfying (1). Does this imply  $|\mathcal{F}| \leq \binom{n}{s}$ ?

**PROBLEM 2.** Suppose  $\mathcal{F} = \{F_1, \dots, F_f\}$  satisfies  $|F_i - F_j| \leq s$  for  $1 \leq i < j \leq f$ . What is  $\max f$ ? If we assume that  $\mathcal{F}$  is an antichain then the proof of Theorem 1 shows that (2) holds. However, in general one can construct much larger families: take all subsets of  $X$  whose size is either at most  $s$  or at least  $n - s$ .

Fix a chain  $E_1 \subset E_2 \subset \dots \subset E_{n-2s-1} \subset X$ ,  $|E_i| = i$ , and add to the preceding all  $(s+i)$ -subsets of  $X$ , containing  $E_i$ ,  $i = 1, \dots, n-2s-1$ . Finally order these sets in increasing order of cardinality. With Z. Füredi we conjecture that this construction is optimal.

ADDED IN PROOF. Z. Füredi, J. Pach and the author solved Problem 2 for  $s=1$  and showed that the above construction is not far from best possible in general, see [5].

#### REFERENCES

- [1] DELSARTE, P., Four fundamental parameters of a code and their combinatorial significance, *Information and Control* **23** (1973), 407—438. *MR* **48** # 13453.
- [2] DELSARTE, P., Bounds for unrestricted codes, by linear programming, *Philips Res. Rep.* **27** (1972), 272—289. *MR* **47** # 3096.
- [3] FRANKL, P. and WILSON, R. M., Intersection theorems with geometric consequences, *Combinatorica* **1** (1981), 357—368.
- [4] KATONA, G. O. H., Open problems, *Matematikus Kurir*, 1983 (in Hungarian).
- [5] FRANKL, P., FÜREDI, Z. and PACH, J., Bounding one-sided differences, *Graphs and combinatorics* (submitted).

(Received August 8, 1983)

E. R. 175 COMBINATOIRE  
CENTRE DE MATHÉMATIQUE SOCIALE  
54, BOULEVARD RASPAIL  
F—75270 PARIS Cedex 06  
FRANCE

# СТРОГО НАСЛЕДСТВЕННЫЕ РАДИКАЛЫ В КЛАССЕ ВСЕХ ТОПОЛОГИЧЕСКИХ КОЛЕЦ

TRINH ĐĂNG KHÔI

В настоящей работе изучаются строго наследственные радикалы, т. е., радикалы, для которых подкольцо  $A$  радикального кольца  $R$  радикально. Известно, что в классе всех топологических колец любой наследственный, и значит строго наследственный радикал либо подыдемпотентен, либо наднильпотентен. Подыдемпотентный наследственный радикал полностью описан в [1]. Нами доказано, что в классе всех топологических колец нет нетривиальных наднильпотентных строго наследственных радикалов (теорема 2).

Мы в качестве класса  $K$  рассматриваем класс всех топологических колец.

Пусть  $R$ —кольцо с единицей. Через  $R[x]$  обозначим кольцо многочленов от коммутирующих с элементами из  $R$  и между собой счетного множества  $X$  переменных с коэффициентами из  $R$ ,  $R_0[X]$  подкольцо кольца  $R[X]$ , состоящее из многочленов без свободных членов. Определим

$$\mathcal{T} = \sum_{i=1}^{\infty} R[X]x_i^{2^i}$$

$$S_n = R[x_1, x_2, \dots, x_{n-1}]$$

и

$$\mathcal{V}_n = \sum_{k=n}^{\infty} \sum_{i=1}^{[k/n]} S_k(x_k^i - 1) + \mathcal{T}$$

Отметим некоторые свойства множеств  $\mathcal{V}_n$ .

$$1. \quad \bigcap_{n=1}^{\infty} \mathcal{V}_n = \mathcal{T}.$$

Ясно, что  $\bigcap_{n=1}^{\infty} \mathcal{V}_n \supseteq \mathcal{T}$ .

Пусть  $\varphi \notin \mathcal{T}$  и  $m$ —такое число, что  $\varphi$  не зависит от  $x_i$ ,  $i \geq m$ . Покажем, что  $\varphi \notin \mathcal{V}_m$ . Допустим  $\varphi \in \mathcal{V}_m$ . Тогда

$$\varphi = \sum_{k=m}^r \sum_{i=1}^{[k/m]} f_{k,i}(x_k^i - 1) + \psi$$

где  $f_{k,i} \in S_k$ ,  $\psi \in \mathcal{T}$ . Можно считать, что  $f_{k,i} = \sum_{j=1}^p U_{k,i,j}$  причём каждый од-

ночлен  $U_{k,i,j} \notin \mathcal{F}$ , и значит в каждом из  $U_{k,i,j}$  переменные  $x_l$  входят в степени меньше  $2^l$  и  $U_{k,i,j}$  не зависит от  $x_l$ ,  $l \geq k$ . Тогда ни один из одночленов, входящих в  $\sum_{k=m}^r \sum_{i=1}^{[k/m]} f_{k,i}(x_k^i - 1)$  не может уничтожиться ни с каким одночленом, входящим в  $\psi$ . Так как  $\varphi$  не зависит от  $x_i$  для  $i \geq m$ , то

$$\sum_{k=m}^r \sum_{i=1}^{[k/m]} f_{k,i}(x_k^i - 1) = 0$$

и значит  $\varphi = \psi \in \mathcal{F}$ .

Получили противоречие. Следовательно,  $\varphi \notin \mathcal{V}_m$ . Из произвольности  $\varphi$  следует, что  $\bigcap_{n=1}^{\infty} \mathcal{V}_n \subseteq \mathcal{F}$  и значит  $\bigcap_{n=1}^{\infty} \mathcal{V}_n = \mathcal{F}$ .

2. Для любых множеств  $\mathcal{V}_m, \mathcal{V}_n$  существует такое  $\mathcal{V}_p$ , что  $\mathcal{V}_p \subseteq \mathcal{V}_m \cap \mathcal{V}_n$ .

В самом деле, выбираем  $p \geq \max(m, n)$ . Ясно, что  $\mathcal{V}_p \subseteq \mathcal{V}_m \cap \mathcal{V}_n$ . Легко проверить следующие:

3.  $\mathcal{V}_n + \mathcal{V}_n \subseteq \mathcal{V}_n$  и  $-\mathcal{V}_n \subseteq \mathcal{V}_n$  для любого множества  $\mathcal{V}_n$ .

4.  $\mathcal{V}_{2m} \cdot \mathcal{V}_{2m} \subseteq \mathcal{V}_m$ .

В самом деле, пусть  $\varphi = \varphi_1 \cdot \varphi_2 \in \mathcal{V}_{2m} \cdot \mathcal{V}_{2m}$ , где  $\varphi_i \in \mathcal{V}_{2m}$ ,  $i=1, 2$ .

$$\begin{aligned} \varphi &= \left( \sum_{k=2m}^{r_1} \sum_{i=1}^{[k/2m]} U_{k,i,j}(x_k^i - 1) + \psi_1 \right) \left( \sum_{t=2m}^{r_2} \sum_{p=1}^{[t/2m]} W_{t,p,h}(x_t^p - 1) + \psi_2 \right) = \\ &= \sum_{k,i,j,t,p,h} (U_{k,i,j} W_{t,p,h}(x_k^i - 1)(x_t^p - 1) + U_{k,i,j}(x_k^i - 1)\psi_2 + \psi_1 W_{t,p,h}(x_t^p - 1) + \psi_1 \psi_2) = \\ &= \sum_{k,i,j,t,p,h} (U_{k,i,j} W_{t,p,h}(x_k^i x_t^p - x_k^i - x_t^p + 1) + U_{k,i,j}(x_k^i - 1)\psi_2 + \psi_1 W_{t,p,h}(x_t^p - 1) + \psi_1 \psi_2). \end{aligned}$$

Ясно, что

$$U_{k,i,j}(x_k^i - 1)\psi_2 + \psi_1 W_{t,p,h}(x_t^p - 1) + \psi_1 \psi_2 \in \mathcal{F} \subseteq \mathcal{V}_m$$

для всех  $k, i, j, t, p, h$ . Для доказательства того, что  $\varphi \in \mathcal{V}_m$  достаточно доказать, что

$$U_{k,i,j} W_{t,p,h}(x_k^i x_t^p - x_k^i - x_t^p + 1) \in \mathcal{V}_m$$

для всех  $k, i, j, t, p, h$ .

Рассмотрим два подслучая.

а) Если  $k=t$ , то

$$\begin{aligned} &U_{k,i,j} W_{t,p,h}(x_k^i x_t^p - x_k^i - x_t^p + 1) = \\ &= U_{k,i,j} W_{k,p,h}(x_k^{i+p} - x_k^i - x_k^p + 1) = U_{k,i,j} W_{k,p,h}((x_k^{i+p} - 1) - (x_k^i - 1) - (x_k^p - 1)). \end{aligned}$$

Так как  $i \leq \left\lfloor \frac{k}{2m} \right\rfloor$ ,  $p \leq \left\lfloor \frac{k}{2m} \right\rfloor$ , то  $i+p \leq 2 \left\lfloor \frac{k}{2m} \right\rfloor \leq \left\lfloor \frac{k}{m} \right\rfloor$  и  $i \leq \left\lfloor \frac{k}{2m} \right\rfloor$ ,  $p \leq \left\lfloor \frac{k}{2m} \right\rfloor$ . Следовательно,

$$U_{k,i,j} W_{k,p,h}((x_k^{i+p} - 1) - (x_k^i - 1) - (x_k^p - 1)) \in \mathcal{V}_m$$

и значит

$$U_{k,i,j} W_{t,p,h}(x_k^i x_t^p - x_k^i - x_t^p + 1) \in \mathcal{V}_m.$$

б) Если  $k \neq i$ , допустим для определенности, что  $k < i$ . Тогда

$$U_{k,i,j}W_{i,p,h}(x_k^i x_i^p - x_k^i - x_i^p + 1) = U_{k,i,j}W_{i,p,h}x_k^i(x_i^p - 1) - U_{k,i,j}W_{i,p,h}(x_i^p - 1).$$

Так как  $p \leq \left\lfloor \frac{i}{2m} \right\rfloor < \left\lfloor \frac{i}{m} \right\rfloor$ , то

$$U_{k,i,j}W_{i,p,h}x_k^i(x_i^p - 1) - U_{k,i,j}W_{i,p,h}(x_i^p - 1) \in \mathcal{V}_m$$

и значит

$$U_{k,i,j}W_{i,p,h}(x_k^i x_i^p - x_k^i - x_i^p + 1) \in \mathcal{V}_m.$$

Следовательно,  $\varphi \in \mathcal{V}_m$ , что и требовалось.

5. Для любого множества  $\mathcal{V}_n$  и произвольного элемента  $f \in R[X]$  существует такое  $\mathcal{V}_m$ , что  $f\mathcal{V}_m \subseteq \mathcal{V}_n$  и  $\mathcal{V}_m f \subseteq \mathcal{V}_n$ .

Пусть  $f = \sum_{t=1}^d W_t$ , где  $W_t$  — одночлен, и  $h$  — такое натуральное число, что  $f$  не зависит от  $x_i$ , для  $i \geq h$ . Выбираем  $m > \max(h, n)$ . Докажем теперь, что  $f\mathcal{V}_m \subseteq \mathcal{V}_n$ . Пусть  $\varphi$  — произвольный элемент в  $\mathcal{V}_m$ :

$$\varphi = \sum_{k=m}^r \sum_{i=1}^{|k/m|} \sum_{j=1}^p U_{k,i,j}(x_k^i - 1) + \psi,$$

где  $\psi \in \mathcal{T}$ . Тогда

$$f\varphi = \sum_{k=m}^r \sum_{i=1}^{|k/m|} \sum_{j=1}^p \sum_{t=1}^d (U_{k,i,j}W_t(x_k^i - 1) + \psi W_t) \in \mathcal{V}_n,$$

ибо  $\psi W_t \in \mathcal{T}$ , а  $U_{k,i,j}W_t$  не зависит от  $x_i$  для  $i \geq k$ .

Второе включение аналогично доказывается.

6.  $\mathcal{V}_1 \cap (R + \mathcal{T}) = \mathcal{T}$ .

Ясно, что  $\mathcal{T} \subseteq \mathcal{V}_1 \cap (R + \mathcal{T})$ . Пусть

$$\varphi \in \mathcal{V}_1 \cap (R + \mathcal{T}), \quad \varphi = \sum_{k=1}^r \sum_{i=1}^k f_{k,i}(x_k^i - 1) + \psi,$$

где  $\psi \in \mathcal{T}$ . Можем считать, что  $f_{k,i} = \sum_{j=1}^p U_{k,i,j}$ , причем каждый одночлен  $U_{k,i,j} \notin \mathcal{T}$ , и значит в каждом из  $U_{k,i,j}$  переменные  $x_i$  входят в степени меньше, чем  $2^i$  и  $U_{k,i,j}$  не зависит от  $x_n$ , для  $n \geq k$ . Тогда

$$\varphi = \sum_{k=1}^r \sum_{i=1}^k \sum_{j=1}^p U_{k,i,j}(x_k^i - 1) + \psi.$$

Многочлен

$$\sum_{k=1}^r \sum_{i=1}^k \sum_{j=1}^p U_{k,i,j}(x_k^i - 1)$$

либо равен нулю, либо является суммой одночленов, в которых каждый из переменных  $x_i$  входит в степени меньше, чем  $2^i$  и значит каждый из этих одно-

членов не принадлежит  $\mathcal{T}$ . Так как  $\varphi \in R + \mathcal{T}$ , то

$$\sum_{k=1}^r \sum_{i=1}^k \sum_{j=1}^p U_{k,i,j} (x_k^i - 1) \in R.$$

Поскольку элементы из  $R$  не зависят от  $x_k$ , то

$$\sum_{k=1}^r \sum_{i=1}^k \sum_{j=1}^p U_{k,i,j} (x_k^i - 1) = 0,$$

и значит  $\varphi \in \mathcal{T}$ . Из произвольности  $\varphi$  следует, что  $\mathcal{V}_1 \cap (R + \mathcal{T}) \subseteq \mathcal{T}$ . Равенство доказано.

Рассмотрим теперь фактор-кольца

$$\overline{R[X]} = R[X]/\mathcal{T}, \quad \overline{R_0[X]} = R_0[X]/\mathcal{T}, \quad \overline{R} = (R + \mathcal{T})/\mathcal{T}$$

и совокупность

$$\{\overline{\mathcal{V}_n} = \mathcal{V}_n/\mathcal{T} \mid n = 1, 2, \dots\}.$$

**Теорема 1.** Совокупность множеств  $\{\overline{\mathcal{V}_n}\}_{n=1,2,\dots}$  можно взять в качестве базиса окрестностей нуля в кольце  $\overline{R[X]}$ , чтобы превратить  $\overline{R[X]}$  в топологическое кольцо. При этом  $\overline{R}$  является дискретным подкольцом в  $\overline{R[X]}$ , а  $\overline{R_0[X]}$  является всюду плотным подкольцом в  $\overline{R[X]}$ .

**Доказательство.** Из свойств 1—6 следует справедливость двух первых утверждений. Остается доказать, что  $\overline{R_0[X]}$ —всюду плотное подкольцо в  $\overline{R[X]}$ .

Пусть  $\overline{\varphi}$ —произвольный элемент в  $\overline{R[X]}$ , и  $\overline{\mathcal{V}_n}$ —произвольная базисная окрестность нуля в  $\overline{R[X]}$ . Докажем теперь, что

$$(\overline{\varphi} + \overline{\mathcal{V}_n}) \cap \overline{R_0[X]} \neq \emptyset.$$

В самом деле, пусть  $\varphi$ —некоторый прообраз элемента  $\overline{\varphi}$  в  $R[X]$ , причем  $\varphi = \varphi_0 + a$ , где  $a \in R$ ,  $\varphi_0$ —многочлен без свободного члена. Тогда

$$\varphi = \varphi_0 + a = \varphi_0 + ax_n - a(x_n - 1).$$

Так как  $a(x_n - 1) \in \mathcal{V}_n$ , то

$$\varphi_0 + ax_n \in R_0[X] \cap (\varphi + \mathcal{V}_n)$$

и значит

$$(\overline{\varphi} + \overline{\mathcal{V}_n}) \cap \overline{R_0[X]} \neq \emptyset.$$

Из произвольности элемента  $\overline{\varphi}$  и окрестности  $\overline{\mathcal{V}_n}$  следует, что  $\overline{R_0[X]}$ —всюду плотный идеал в  $\overline{R[X]}$ . Этим доказательство закончивается.

**Теорема 2.** Если  $\mathcal{L}$ —надильпотентный строго наследственный радикал в классе  $K$  всех топологических колец, то всякое топологическое кольцо из  $K$  является  $\mathcal{L}$ -радикальным кольцом.

**Доказательство.** Пусть  $R'$ —любое кольцо из  $K$  и  $R$ —такое дискретное кольцо с единицей, что  $R'$  в дискретной топологии является подкольцом в  $R$ .



Так как  $\mathcal{L}$  — строго наследственный радикал, то достаточно доказать, что  $R$  является  $\mathcal{L}$ -радикальным кольцом.

Пусть  $R[X]$  — как и раньше — кольцо многочленов от коммутирующих с элементами из  $R$  и между собой счетного множества  $X$  переменных  $x_1, x_2, \dots$  с коэффициентами из  $R$ ,  $R_0[X]$  — подкольцо кольца  $R[X]$ , состоящее из многочленов без свободных членов. Определим

$$R_n[X] = R[X]x_n \quad \text{и} \quad \mathcal{T} = \sum_{i=1}^{\infty} R[X]x_i^{2^i}.$$

Так как  $x_i^{2^i} \in \mathcal{T}$ , то  $(R_n[X])^{2^n} \subseteq \mathcal{T}$  и значит  $(R_n[X] + \mathcal{T})/\mathcal{T}$  является нильпотентным идеалом в  $\overline{R[X]} = R[X]/\mathcal{T}$ .

В силу наднильпотентности радикала  $\mathcal{L}$  следует, что  $(R_n[X] + \mathcal{T})/\mathcal{T}$  является  $\mathcal{L}$ -радикальным кольцом.

Тогда  $\sum_{n=1}^{\infty} (\mathcal{T} + R_n[X])/\mathcal{T}$  будет  $\mathcal{L}$ -радикальным кольцом.

Так как

$$\sum_{n=1}^{\infty} (R_n[X] + \mathcal{T})/\mathcal{T} = \left( \sum_{n=1}^{\infty} R_n[X] + \mathcal{T} \right) / \mathcal{T} = R_0[X]/\mathcal{T}$$

и в силу теоремы 1,  $R_0[X]/\mathcal{T}$  — всюду плотное подкольцо в  $R[X]/\mathcal{T}$ , то  $R[X]/\mathcal{T}$  будет  $\mathcal{L}$ -радикальным кольцом. Тогда  $R$  как подкольцо кольца  $R[X]/\mathcal{T}$ , будет  $\mathcal{L}$ -радикальным кольцом. Этим теорема полностью доказана.

Следствие. Если  $\mathcal{L}$ -нетривиальный строго наследственный радикал в классе всех топологических колец, то существует такой нетривиальный алгебраический радикал  $\mathcal{L}^*$ , что  $\mathcal{L} \subseteq \mathcal{L}^*$ .

Доказательство. Так как  $\mathcal{L}$  — строго наследственный радикал, то либо  $\mathcal{L}$  — наднильпотентный строго наследственный радикал либо  $\mathcal{L}$  — подидемпотентный строго наследственный радикал (см. [1] теорема 8).

Если  $\mathcal{L}$  — наднильпотентный строго наследственный радикал, то сам  $\mathcal{L}$  является тривиальным алгебраическим радикалом, что противоречит нетривиальности радикала  $\mathcal{L}$ .

Если же  $\mathcal{L}$  — подидемпотентный строго наследственный радикал, то существуют такие наборы простых чисел  $p_1, \dots, p_k$  и натуральных чисел  $n_1, \dots, n_k$ , что все  $\mathcal{L}$ -радикальные кольца содержатся в классе  $T(p_1^{n_1}, p_2^{n_2}, \dots, p_k^{n_k})$  (см. [1], теорема 8).

Известно, что класс  $T(p_1^{n_1}, \dots, p_k^{n_k})$  является классом всех  $\mathcal{L}^*$ -радикальных колец для некоторого алгебраического радикала  $\mathcal{L}^*$  (см. [2], теорема 3). Тогда  $\mathcal{L} \subseteq \mathcal{L}^*$ , что и требовалось.

#### ЛИТЕРАТУРА

- [1] ARNAUTOV, V. I. and VODINČAR, M. I., Radicals of topological rings, *Mat. Issled.* **3** (1968), вып. 2 (8), 31—61 (in Russian). *MR* **40** # 7313.
- [2] ARNAUTOV, V. I., Algebraic radicals in topological rings (supplement). Topological structures and algebraic systems, *Mat. Issled.*, Вып. **44** (1977), 36—41, 178. *MR* **58** # 28051.

(Поступило 10-ого августа 1983 г.)

FACULTY OF MATHEMATICS  
HANOI PEDAGOGICAL INSTITUTE  
HANOI  
VIETNAM



# CONGRUENCE UNIFORM ALGEBRAS WITH PSEUDOCOMPLEMENTATION

R. BEAZER

## 1. Introduction

In [13], W. Taylor initiated the study of varieties of algebras with uniform congruences and various related notions together with their connection with congruence regularity. A congruence relation on an algebra is uniform if all its congruence classes have the same cardinality and an algebra is called congruence uniform if every one of its congruences is uniform. In [13], it was stated, without proof, that every member of any variety generated by a quasi-primal algebra is congruence uniform. R. McKenzie [10] showed that every finite member of any directly representable variety is congruence uniform. Apart from these universal algebraic results very little seems to be in literature about congruence uniform algebras in common or garden varieties. Of course, every Boolean algebra is congruence uniform. The aim of this note is to partially fill that gap for the various varieties of (semi-) lattices with pseudocomplementation. The congruence uniform algebras in the varieties of pseudocomplemented semilattices, lattices with pseudocomplementation ( $p$ -algebras) and bounded relatively pseudocomplemented lattices (Heyting algebras) are easily identified as the Boolean algebras. Put another way, an algebra in any of the aforementioned varieties is congruence uniform if and only if it is congruence regular. The main result in this note shows that exactly the same is true for double  $p$ -algebras. It is further shown that every finite double Heyting algebra is congruence uniform.

## 2. Preliminaries

A *pseudocomplemented semilattice* ( $p$ -semilattice) is an algebra  $\langle L, \wedge, *, 0 \rangle$ , where  $\langle L, \wedge, 0 \rangle$  is a semilattice with least element 0 and  $*$  is a unary operation, called *pseudocomplementation*, defined on  $L$  by

$$a \wedge x = 0 \Leftrightarrow x \leq a^*.$$

A *lattice with pseudocomplementation* ( $p$ -algebra) is an algebra  $\langle L, \vee, \wedge, *, 0, 1 \rangle$ , where  $\langle L, \vee, \wedge, 0, 1 \rangle$  is a bounded lattice and  $*$  is pseudocomplementation.

A *Heyting algebra* is an algebra  $\langle L, \vee, \wedge, *, 0, 1 \rangle$ , where  $\langle L, \vee, \wedge, 0, 1 \rangle$  is a bounded lattice and  $*$  is a binary operation, called *implication* (or *relative pseudo-*

---

1980 *Mathematics Subject Classification*. Primary 06D15; Secondary 06D20.

*Key words and phrases*. Double  $p$ -algebra, double Heyting algebra, congruence uniform, congruence regular.

complementation), defined on  $L$  by

$$a \wedge x \leq b \leftrightarrow x \leq a * b.$$

The subset  $B(L) = \{x \in L; x = x^{**}\}$  of a  $p$ -semilattice  $\langle L, \wedge, *, 0 \rangle$  is a Boolean algebra under the operations  $\wedge, \cup, *$  where the join operation  $\cup$  is defined on  $B(L)$  by  $a \cup b = (a^* \wedge b^*)^*$ . In the event that  $L$  is a  $p$ -algebra,  $a \cup b = (a \vee b)^{**}$ , for any  $a, b \in B(L)$ . The relation  $\varphi$  defined on any of the aforementioned algebras by

$$a \equiv b(\varphi) \leftrightarrow a^* = b^*$$

is a congruence, called the *Glivenko congruence*, and  $L/\varphi \cong B(L)$ .

A *double  $p$ -algebra* is an algebra  $\langle L, \vee, \wedge, *, +, 0, 1 \rangle$  in which the unary operation  $+$  is characterized by

$$a \vee x = 1 \leftrightarrow x \equiv a^+$$

and whose deletion yields a  $p$ -algebra. The relation  $\Phi$  defined on any double  $p$ -algebra by

$$a \equiv b(\Phi) \leftrightarrow a^* = b^* \quad \text{and} \quad a^+ = b^+$$

is a congruence, called the *determination congruence*.

A *double Heyting algebra* is an algebra  $\langle L, \vee, \wedge, *, +, 0, 1 \rangle$  in which the binary operation  $+$  is characterized by

$$a \vee x \leq b \leftrightarrow x \leq a + b$$

and whose deletion yields a Heyting algebra.

For all unexplained lattice theoretic notation and terminology we refer to [6]. For the standard rules of computation in  $p$ -semilattices,  $p$ -algebras, Heyting algebras and double  $p$ -algebras we refer to [1], [6] and [10].

If  $\theta$  is a congruence on an algebra  $A$  and  $B$  is a non-empty subset of  $A$ , we write  $[a]_\theta$  for the congruence class  $\{x \in A; x \equiv a(\theta)\}$  and  $\Theta(B)$  for the smallest congruence on  $A$  collapsing  $B$ . The congruence lattice of  $A$  will be denoted by  $\text{Con}(A)$  and  $\omega, \iota$  will denote its least and greatest elements respectively. For all undefined universal algebraic terminology and background we refer to [5] or [7].

### 3. Uniformity and regularity

We start by recalling the definition of regularity. An algebra is *regular* if any two of its congruences are equal as soon as they have a congruence class in common. Now, if  $L$  is a  $p$ -semilattice,  $p$ -algebra or Heyting algebra and  $\varphi$  is its Glivenko congruence then  $[0]_\varphi = \{0\}$  and  $\varphi = \omega$  if and only if  $L$  is Boolean. This, in conjunction with the fact that  $p$ -semilattice,  $p$ -algebra and Heyting algebra congruences on a Boolean algebra are Boolean algebra congruences, immediately yields the following:

**THEOREM 1.** *For a  $p$ -semilattice,  $p$ -algebra or Heyting algebra  $L$  the following are equivalent:*

- (1)  $L$  is congruence uniform,
- (2)  $L$  is regular,
- (3)  $\varphi = \omega$ ,
- (4)  $L$  is Boolean.

For the rest of this section we will be concerned with double  $p$ -algebras. J. Varlet [14] showed that for a double  $p$ -algebra to be regular it is necessary and sufficient that its determination congruence  $\Phi = \omega$ . In [8], T. Katriňák showed that every regular double  $p$ -algebra is a double Heyting algebra (and, therefore, distributive) in which the operation  $*$  is given by

$$x * y = (x^* \vee y^{**})^{**} \wedge [(x \vee x^*)^+ \vee x^* \vee y \vee y^*]$$

and  $x + y$  is given by the dual expression. Various other useful characterizations of regular double  $p$ -algebras may be found in [3], [8] and [9].

**LEMMA 2.** *Let  $L$  be a Heyting algebra and let  $x, a, c \in L$ . Then the following are equivalent:*

- (1)  $x \wedge c = a \wedge c$ ,
- (2)  $(x * a) \wedge (a * x) \cong c$ ,
- (3)  $x \in [a \wedge c, c * a]$ .

*If, in addition,  $c$  has a complement  $c'$  in  $L$  then  $c * a = c' \vee a$ .*

**PROOF.** The equivalence of (1) and (2) is well-known (see [11], for example) while the equivalence of (1) and (3) together with the last part follow in a straightforward manner from the definition of  $*$ .

A normal filter in double  $p$ -algebra is a lattice filter  $F$  having the property that  $f^{**} \in F$  whenever  $f \in F$ .

**LEMMA 3.** *If  $F$  is a normal filter in a distributive double  $p$ -algebra  $L$  then  $\Theta(F) = \Theta_{\text{lat}}(F)$ , the smallest lattice congruence on  $L$  collapsing  $F$ , and  $F = [1] \Theta(F)$ .*

**PROOF.** It is well-known that, for any filter  $F$  in a distributive lattice  $L$ ,  $a \equiv b(\Theta_{\text{lat}}(F)) \Leftrightarrow a \wedge f = b \wedge f$ , for some  $f \in F$ , and that  $\Theta_{\text{lat}}(F)$  preserves  $*$  in the event that  $L$  is a distributive  $p$ -algebra (see [1], for example). Therefore, it is enough to show that  $\Theta_{\text{lat}}(F)$  preserves  $+$ . To effect this, suppose that  $f \in F$  and  $a \wedge f = b \wedge f$  then  $a^+ \vee f^+ = (a \wedge f)^+ = (b \wedge f)^+ = b^+ \vee f^+$  from which it follows, on taking the meet of both sides with  $f^{**}$  and using distributivity, that  $a^+ \wedge f^{**} = b^+ \wedge f^{**}$ . Therefore  $a^+ \equiv b^+(\Theta_{\text{lat}}(F))$ , since  $f^{**} \in F$ . The last part is an immediate consequence of the aforementioned description of  $\Theta_{\text{lat}}(F)$ .

**THEOREM 4.** *A double  $p$ -algebra is congruence uniform if and only if it is congruence regular.*

**PROOF.** In one direction the proof is trivial. If  $L$  is a congruence uniform double  $p$ -algebra,  $\Phi$  is the determination congruence on  $L$  and  $a \in L$  then  $[a] \Phi = [0] \Phi = |\{0\}| = 1$  so that  $\Phi = \omega$  and, therefore,  $L$  is regular.

Suppose, now, that  $L$  is a regular double  $p$ -algebra. Then  $L$  is a double Heyting algebra in which  $*$  and  $+$  are double  $p$ -algebra polynomials. For  $x, y \in L$  let us write

$$x \circ y = (x * y) \wedge (y * x), \quad x \oplus y = (x + y) \vee (y + x)$$

and observe that

$$x \circ x = 1, \quad 1 \circ x = x, \quad x \oplus x = 0, \quad 0 \oplus x = x.$$

Define double  $p$ -algebra polynomials  $f^\circ, f^\oplus: L^3 \rightarrow L$  and  $f: L^3 \rightarrow L^2$  by

$$f^\circ(x, y, z) = (x \circ y) \circ z, \quad f^\oplus(x, y, z) = (x \oplus y) \oplus z$$

and

$$f(x, y, z) = \langle f^\circ(x, y, z), f^\oplus(x, y, z) \rangle.$$

Let  $\theta \in \text{Con}(L)$  and let us suppose, in the first instance, that no congruence class of  $\theta$  is finite. We claim that  $x \mapsto f(x, a, b)$  is a one-to-one map from  $[a]\theta$  to  $([b]\theta)^2$ , for any  $a, b \in L$ . Indeed, if  $x \in [a]\theta$  then

$$\begin{aligned} f^\circ(x, a, b) &\equiv f^\circ(a, a, b)(\theta) \\ &= (a \circ a) \circ b = 1 \circ b = b \end{aligned}$$

and so  $f^\circ(x, a, b) \in [b]\theta$ . Similarly,  $f^\oplus(x, a, b) \in [b]\theta$ . Therefore,  $x \mapsto f(x, a, b)$  maps  $[a]\theta$  to  $([b]\theta)^2$ . Furthermore, if  $x, y \in [a]\theta$  and  $f(x, a, b) = f(y, a, b)$  then  $(f^\circ(x, a, b))^{**} = (f^\circ(y, a, b))^{**}$ . But  $x \mapsto x^{**}$  is a Heyting algebra homomorphism from  $L$  onto  $B(L)$  and so  $(f^\circ(x, a, b))^{**} = (x^{**} \circ_B a^{**}) \circ_B b^{**}$ , where  $z \circ_B w = (z^* \cup w) \wedge (w^* \cup z)$  for any  $z, w \in B(L)$ . It follows, now, that  $x^{**} = y^{**}$  because  $\langle B(L), \circ_B, 1 \rangle$  is a group satisfying the identity  $z \circ_B z = 1$ . By dual reasoning, we also have  $x^{++} = y^{++}$  and so  $x = y$ , since  $\Phi = \omega$ . Thus,  $|[a]\theta| \leq |[b]\theta|^2 = |[b]\theta|$ , since  $[b]\theta$  is infinite. Similarly,  $|[b]\theta| \leq |[a]\theta|$  and so  $|[a]\theta| = |[b]\theta|$ . Next, let us suppose that  $\theta$  has a finite congruence class, say  $[b]\theta$ . Then every congruence class of  $\theta$  must be finite, since by the preceding argument  $|[a]\theta| \leq |[b]\theta|^2$ , for any  $a \in L$ . In particular,  $[1]\theta$  is finite and so  $[1]\theta$  is a principal filter generated by some complemented element  $c$  in  $L$ . Indeed, it is obvious that  $[1]\theta = [c]$  for some  $c \in L$ , and that  $c^{++} \equiv c$ , since  $[1]\theta$  is a normal filter in  $L$ . Therefore,  $c \wedge c^+ \leq c^{++} \wedge c^+ = 0$  from which it follows that  $c^+$  is the complement of  $c$  in  $L$ . By Lemma 3, we have

$$[1]\theta([c]) = [c] = [1]\theta$$

and so  $\theta = \theta([c])$ , by regularity. As a simple consequence of this and Lemmas 2 and 3, we have  $[a]\theta = [a \wedge c, c' \vee a]$ , for any  $a \in L$ . We claim that, for any  $a \in L$ ,  $x \mapsto (c' \vee a) \wedge x$  is a one-to-one map from  $[c]$  onto  $[a]\theta$ . Clearly,  $(c' \vee a) \wedge x \in [a \wedge c, c' \vee a]$ , whenever  $x \equiv c$ . Also, if  $z \in [a \wedge c, c' \vee a]$  and  $x = z \vee c$  then  $x \equiv c$  and, by distributivity,

$$(c' \vee a) \wedge x = [(c' \vee a) \wedge z] \vee [(c' \vee a) \wedge c] = z \vee (a \wedge c) = z.$$

Therefore,  $x \mapsto (c' \vee a) \wedge x$  is onto  $[a \wedge c, c' \vee a]$ . Moreover, if  $x, y \equiv c$  and  $(c' \vee a) \wedge x = (c' \vee a) \wedge y$  then  $x \circ y \equiv c' \vee a$ , by Lemma 2, and  $x \circ y = (x * y) \wedge (y * x) \equiv y \wedge x \equiv c$  so that  $x \circ y \equiv (c' \vee a) \vee c = 1$ . Therefore  $x = y$  and we conclude that  $x \mapsto (c' \vee a) \wedge x$  is one-to-one. It follows, now, that  $|[1]\theta| = |[a]\theta|$ , for any  $a \in L$ . Thus,  $L$  is congruence uniform.

**COROLLARY 5.** *Any Boolean algebra is congruence uniform.*

The proof of Theorem 4 breaks down for an arbitrary double Heyting algebra owing to the failure of the determination principle:  $x^* = y^*$  and  $x^+ = y^+ \Rightarrow x = y$ . Nevertheless, we have

**THEOREM 6.** *Any finite double Heyting algebra is congruence uniform.*



PROOF. Let  $L$  be a finite double Heyting algebra. By the subdirect product theorem,  $L$  is a subdirect product of finitely many finite subdirectly irreducible double Heyting algebras. However, according to [4], every finite subdirectly irreducible double Heyting algebra is simple. Moreover, the congruences on any double Heyting algebra permute, since it is well-known that Heyting algebra congruences permute. Therefore, by a standard universal algebraic result (see [5], for example),  $L$  is a direct product of finitely many simple double Heyting algebras  $S_1, S_2, \dots, S_n$ , say. Since the variety of double Heyting algebras is congruence distributive, it follows, by a theorem of H. Werner (see [5]) that every congruence  $\theta$  on  $L$  is a product congruence; that is,  $\theta$  is of the form  $\prod_{i=1}^n \theta_i$  where  $\theta_i \in \text{Con}(S_i)$ ,  $1 \leq i \leq n$ . Thus, if

$\mathbf{a} = \langle a_1, \dots, a_n \rangle$ ,  $\mathbf{b} = \langle b_1, \dots, b_n \rangle \in \prod_{i=1}^n S_i$  then

$$|[\mathbf{a}]\theta| = \left| \prod_{i=1}^n [a_i]\theta_i \right| = \prod_{i=1}^n |[a_i]\theta_i| = \prod_{i=1}^n |[b_i]\theta_i| = |[\mathbf{b}]\theta|,$$

since simple algebras are obviously congruence uniform. This completes the proof that  $L$  is congruence uniform.

#### 4. Concluding remarks

(1) There is only one (non-trivial) finite semilattice that is congruence uniform, namely the two-element semilattice. Indeed, if  $L$  is a finite semilattice,  $|L| > 2$  and  $a$  is an atom of  $L$ , then the equivalence relation  $\theta_a$  induced by the partition  $\{\{0, a\}, \{x\} : x \in L \setminus \{0, a\}\}$  of  $L$  is obviously a non-uniform congruence relation on  $L$ .

(2) A finite distributive lattice is congruence uniform if and only if it is Boolean. This can be argued as in the pre-amble to Theorem 1 because every finite distributive lattice is pseudocomplemented and all lattice congruences on a Boolean lattice preserve complementation. Alternatively, we can argue as follows. Suppose that  $L$  is a finite congruence uniform distributive lattice. Let  $P$  be a prime ideal of  $L$  and let  $M$  be a maximal ideal containing  $P$ . The equivalence relation induced by the partition  $\{P, L \setminus P\}$  of  $L$  is a congruence so that  $|P| = |L \setminus P|$  and, therefore,  $|P| = \frac{1}{2} |L|$ . Likewise,  $|M| = \frac{1}{2} |L|$ , since every maximal ideal is prime. Consequently, every prime ideal of  $L$  is maximal and so  $L$  is Boolean by Nachbin's theorem (see [1], for example).

(3) The above result can be generalized to a wider class of distributive lattices, namely the distributive  $*$ -lattices of T. P. Speed [12]. These are distributive lattices  $L$  with least element 0 satisfying the condition that, for all  $a \in L$ , there exists  $a' \in L$  such that  $(a)^{**} = (a')^*$ , where  $*$  denotes pseudocomplementation in the ideal lattice of  $L$ . The reason why we can generalize to this class is that T. P. Speed [12] has shown that for a distributive lattice with 0 to be a distributive  $*$ -lattice it is necessary and sufficient that  $L/\varphi$  is a Boolean lattice, where  $\varphi$  is the congruence defined on  $L$  by  $a \equiv b(\varphi) \Leftrightarrow (a)^* = (b)^*$ .

(4) Algebras that are direct products of finitely many simple algebras in any congruence distributive variety are congruence uniform. In particular, all finite algebras in any semi-simple arithmetical variety, all finite complemented modular lattices and all finite relatively complemented lattices are congruence uniform. In view of this, we leave the reader with the following open problems.

PROBLEM 1. Which (finite) lattices are congruence uniform?

PROBLEM 2. Which infinite double Heyting algebras are congruence uniform?

#### REFERENCES

- [1] BALBES, R. and DWINGER, P., *Distributive lattices*, University of Missouri Press, Columbia Mo. 1974. *MR* 51 # 10185.
- [2] BEAZER, R., The determination congruence on double  $p$ -algebras, *Algebra Universalis* 6 (1976), 121—129. *MR* 54 # 7371.
- [3] BEAZER, R., Regular double  $p$ -algebras with Stone congruence lattices, *Algebra Universalis* 9 (1979), 238—243. *MR* 80j: 06007.
- [4] BEAZER, R., Subdirectly irreducible double Heyting algebras, *Algebra Universalis* 10 (1980), 220—224. *MR* 81d: 06015.
- [5] BURRIS, S. and SANKAPPANAVAR, H. P., *A course in universal algebra*, Graduate Texts in Mathematics, 78 Springer-Verlag, New York—Berlin, 1981. *MR* 83k: 08001.
- [6] GRÄTZER, G., *General lattice theory*, Pure and Applied Mathematics, Vol. 75. Academic Press, New York, 1978; Mathematische Reihe, Band 52, Birkhäuser Verlag, Basel—Stuttgart, 1978; Akademie Verlag, Berlin, 1978. *MR* 80c: 06001a, 06001b.
- [7] GRÄTZER, G., *Universal algebra*, second edition, Springer-Verlag, New York—Heidelberg, 1979. *MR* 80g: 08001.
- [8] KATRÍŇÁK, T., The structure of distributive double  $p$ -algebras. Regularity and congruences, *Algebra Universalis* 3 (1973), 238—246. *MR* 48 # 10924.
- [9] KATRÍŇÁK, T., Subdirectly irreducible distributive double  $p$ -algebras, *Algebra Universalis* 10 (1980), 195—219. *MR* 81g: 06006.
- [10] MCKENZIE, R., Narrowness implies uniformity, *Algebra Universalis* 15 (1982), 67—85. *MR* 83i: 08003.
- [11] RASIOWA, H., *An algebraic approach to non-classical logics*, Studies in Logic and the Foundations of Mathematics, Vol. 78. North-Holland Publishing Co. Amsterdam—London; American Elsevier Publishing Co., Inc., New York, 1974. *MR* 56 # 5285.
- [12] SPEED, T. P., Some remarks on a class of distributive lattices, *J. Austral. Math. Soc.* 9 (1969), 289—296. *MR* 40 # 69.
- [13] TAYLOR, W., Uniformity of congruences, *Algebra Universalis* 4 (1974), 346—360. *MR* 51 # 12658.
- [14] VARLET, J., A regular variety of type  $\langle 2, 2, 1, 1, 0, 0 \rangle$ , *Algebra Universalis* 2 (1972), 218—223. *MR* 48 # 3824.

(Received September 27, 1983)

DEPARTMENT OF MATHEMATICS  
UNIVERSITY GARDENS  
UNIVERSITY OF GLASGOW  
GLASGOW  
G12 8QW  
SCOTLAND

## COMPLETE BASES IN TOPOLOGICAL SPACES

JOSE L. BLASCO

### Introduction

The word "space" will refer to *Tychonoff spaces*. In [2] the author considers a particular class of bases<sup>1</sup>  $\mathcal{D}$  on a space  $X$ , called complete. They are characterized by the relation  $\beta(v(X, \mathcal{D})) = \omega(X, \mathcal{D})^2$  between their associated Wallman spaces. It is known that the family  $Z(X)$  of all zero-sets in the space  $X$  is always a complete base on  $X$ . In some cases each base on  $X$  is complete, for example when the Hewitt realcompactification  $vX$  is Lindelöf ([6], (4.4) and [2], Theorem 5). However, the existence of noncomplete bases is not trivial. For examples of noncomplete bases see [15] and [3]. Other examples have been given in [7], [8], [9] and [11] but as certain inverse-closed subalgebras of  $C(X)$ . It should be noted that all these examples have been constructed on uncountable discrete spaces.

In this paper, a method is provided for constructing noncomplete bases on a wide class of topological spaces. This method is then applied to characterize the Lindelöf property of the Hewitt realcompactification of a space. Also, we shall prove that a paracompact space  $X$  is Lindelöf if and only if each base on  $X$  is complete.

In the last section we give some applications to certain inverse-closed subalgebras of  $C(X)$ .

### Preliminaries

As usual,  $C(X)$  will denote the ring of all continuous real-valued functions on a space  $X$ . We write  $C_I(X)$  for the set  $\{f \in C(X) : f(X) \subseteq [0, 1]\}$ . For any  $f$  and  $g$  in  $C(X)$ , we denote  $(f \vee g)(x) = \max(f(x), g(x))$  and  $(f \wedge g)(x) = \min(f(x), g(x))$ . The set of points of  $X$  where a member  $f$  of  $C(X)$  is equal to zero is called the zero-set of  $f$  and will be denoted by  $Z(f)$ . The collection of all zero-sets of  $X$  will be denoted by  $Z(X)$ . If  $\mathcal{E}$  is a class of subsets of  $X$  and  $A$  is a fixed subset of  $X$ , we shall denote by  $\mathcal{E} \cap A$  the class of all sets of the form  $D \cap A$  with  $D$  in  $\mathcal{E}$ . The  $Q$ -closure of  $A$  in  $X$  is the set  $Q(A, X)$  of all points  $p \in X$  for which each  $G_\delta$ -set about  $p$  meets  $A$ .

Compactifications will always be Hausdorff. We will write  $\beta X$  for the Stone—Čech compactification of the space  $X$ . It is well known that  $vX$  is the  $Q$ -closure of  $X$  in  $\beta X$ .

<sup>1</sup> This notion is due to E. F. Steiner [13] who uses the term separating nest generated intersection ring. An equivalent concept is the strong delta normal base due to R. A. Ald and H. L. Shapiro [1].

<sup>2</sup> Two extensions  $T_1$  and  $T_2$  of a space  $X$  are said to be equivalent if they are homeomorphic via a map that leaves  $X$  pointwise fixed. In this case we write  $T_1 = T_2$ .

1980 *Mathematics Subject Classification*. Primary 54D60; Secondary 54D18.

*Key words and phrases*. Realcompactification, Lindelöf property.

A subset  $B$  of  $X$  is said to be  $z$ -separated from  $A$  if there is a zero-set  $Z$  in  $X$  such that  $B \subset Z$  and  $Z \cap A = \emptyset$ . It is known [12] that  $X$  is a Lindelöf space if and only if there is a compactification  $K$  of  $X$  such that every compact subset of  $K - X$  is  $z$ -separated from  $X$ .

**THEOREM 1.** *Let  $K$  be a compactification of a space  $X$  such that there exist two compact subsets  $C_1$  and  $C_2$  of  $K \sim Q(X, K)$  which are disjoint and homeomorphic. If  $C_2$  is not  $z$ -separated from  $X$ , then there exists a compactification  $K^*$  of  $X$  such that  $K^* \leq K$ ,  $Q(X, K^*) = Q(X, K)$  and  $Z(K^*) \cap X \neq Z(K) \cap X$ .*

**PROOF.** Let  $\phi$  be a homeomorphism from  $C_1$  onto  $C_2$ , let  $K^*$  be the compactification of  $X$  obtained from  $K$  by identifying  $y$  and  $\phi(y)$  for each  $y \in C_1$ , and let  $\psi: K \rightarrow K^*$  be the resulting canonical map. We do not distinguish notationally between  $X$  and  $\psi(X)$ .

First, we shall prove that  $Q(X, K) = Q(X, K^*)$ . Assume that  $p \in K \sim (C_1 \cup C_2 \cup Q(X, K))$ . Then  $p \notin Q(X, K)$ , and so there is  $f \in C_I(K)$  such that  $f(p) = 0$  and  $Z(f) \cap X = \emptyset$ . Let  $g$  be a function in  $C_I(K)$  taking the value 0 on  $p$  and 1 on  $C_1 \cup C_2$ . Since  $\psi$  is quotient, the function  $h$  defined on  $K^*$  by the equality  $h \circ \psi = f \vee g$  is continuous. Moreover  $h(\psi(p)) = 0$  and  $Z(h) \cap X = \emptyset$ , hence  $\psi(p) \notin Q(X, K^*)$ .

Now assume that  $p \in C_1$ . Then  $p \notin Q(X, K)$ , and so there exists  $f_1' \in C_I(K)$  such that  $f_1'(p) = 0$ , while  $f_1'$  vanishes nowhere on  $X$ . Let  $f_1'' \in C_I(K)$  such that  $f_1''(p) = 0$  and  $f_1''(C_2) = 1$ . If  $f_1 = f_1' \vee f_1''$ , then  $f_1(p) = 0$ ,  $Z(f_1) \cap X = \emptyset$  and  $f_1(C_2) = 1$ . Since the set  $C_1 \cup C_2 \cup Z(f_1)$  is compact, there exists  $f_2 \in C_I(K)$  such that  $f_2(C_1 \cup Z(f_1)) = 1$  and  $f_2(\phi(x)) = f_1(x)$  if  $x \in C_1$ . If  $h = f_1 \wedge f_2$ , then  $h \in C(K)$ ,  $h|_{C_1} = f_1|_{C_1}$ ,  $h|_{C_2} = f_2|_{C_2}$  and  $h(p) = h(\phi(p)) = 0$ .

By applying a similar argument to the point  $\phi(p) \in C_2$ , we obtain  $g_1 \in C_I(K)$  such that  $g_1(\phi(p)) = 0$ ,  $Z(g_1) \cap X = \emptyset$  and  $g_1(C_1 \cup Z(f_1)) = 1$ . Let  $g_2 \in C_I(K)$  such that  $g_2(x) = g_1(\phi(x))$  if  $x \in C_1$  and  $g_2(C_2 \cup Z(f_2)) = 1$ . If  $g = g_1 \wedge g_2$ , then  $g|_{C_1} = g_1|_{C_1}$ ,  $g|_{C_2} = g_2|_{C_2}$  and  $g(p) = g(\phi(p)) = 0$ .

The function  $t = h \vee g$  is continuous on  $K$  and  $Z(t) \cap X = \emptyset$  because

$$Z(t) = (Z(f_1) \cup Z(f_2)) \cap (Z(g_1) \cup Z(g_2))$$

and the sets  $Z(f_1) \cap X$ ,  $Z(g_1) \cap X$  and  $Z(g_2) \cap Z(f_2)$  are empty. Therefore, if  $t'$  is the function defined on  $K^*$  by the equality  $t' \circ \psi = t$ , then  $t' \in C(K^*)$ ,  $t'(\psi(p)) = 0$  and  $Z(t') \cap X = \emptyset$ . Hence  $\psi(p) \notin Q(X, K^*)$ . Thus we have proved that  $Q(X, K^*) \subset \psi(Q(X, K))$ . On the other hand, it is easy to check that  $\psi(Q(X, K)) \subset Q(X, K^*)$ , consequently  $\psi(Q(X, K)) = Q(X, K^*)$ . Therefore the restriction of  $\psi$  to  $Q(X, K)$  is a homeomorphism from  $Q(X, K)$  onto  $Q(X, K^*)$  that leaves  $X$  pointwise fixed.

It remains to show that  $Z(K^*) \cap X \neq Z(K) \cap X$ . Let  $Z_1'$  and  $Z_2'$  be disjoint zero-sets in  $K$  such that  $C_1 \subset \text{int}_K Z_1'$  and  $C_2 \subset Z_2'$ . Let us see that  $Z_1' \cap X$  is not the trace of a zero-set in  $K^*$  on  $X$ . In fact, assume there is a zero-set  $Z \in Z(K^*)$  such that  $Z \cap X = Z_1' \cap X$ . Since  $C_1 \subset \text{cl}_K(Z_1' \cap X)$  it follows that  $\psi(C_1) \subset Z$  and consequently  $C_1 \cup C_2 \subset \psi^{-1}(Z)$ . As  $\psi^{-1}(Z) \in Z(K)$  we have that  $Z_2' \cap \psi^{-1}(Z)$  is a zero-set in  $K$  which contains  $C_2$  and misses  $X$ . This is a contradiction because  $C_2$  is not  $z$ -separated from  $X$ . ■

## Complete bases

We will write  $\omega(X, \mathcal{D})$  (resp.  $v(X, \mathcal{D})$ ) for the Wallman compactification (resp. Wallman realcompactification) associated with a given base  $\mathcal{D}$  on  $X$ . For definitions and basic results the reader is referred to [1], [14] and [15]. For a base  $\mathcal{D}$  on  $X$  let  $\hat{\mathcal{D}}$  be the trace on  $X$  of all zero-sets in  $v(X, \mathcal{D})$ . Then  $\hat{\mathcal{D}}$  is a base on  $X$  [13] which contains  $\mathcal{D}$ , since  $\mathcal{D}$  is the trace on  $X$  of all zero-sets in  $\omega(X, \mathcal{D})$ . ([15], Theorem 2.2). A base  $\mathcal{D}$  on  $X$  is called complete if  $\mathcal{D} = \hat{\mathcal{D}}$ . It follows from ([2], Corollary 2.2) that:

- (1)  $\hat{\mathcal{D}}$  is the largest base on  $X$  such that  $v(X, \mathcal{D}) = v(\hat{X}, \mathcal{D})$ .
- (2)  $\hat{\mathcal{D}}$  is the smallest complete base on  $X$  containing  $\mathcal{D}$ .

The following Theorem is the main result.

**THEOREM 2.** *Let  $K$  be a compactification of a space  $X$  such that there exist two compact subsets  $C_1$  and  $C_2$  of  $K \sim Q(X, K)$  which are disjoint and homeomorphic. If  $C_2$  is not  $z$ -separated from  $X$ , then  $Z(K) \cap X$  contains a noncomplete base  $\mathcal{D}$  on  $X$  such that  $\hat{\mathcal{D}} = (Z(K) \cap X)^\wedge$ .*

**PROOF.** According to Theorem 1 there exists a compactification  $K^*$  of  $X$  such that  $K^* \leq K$ ,  $Q(X, K^*) = Q(X, K)$  and  $Z(K) \cap X \neq Z(K^*) \cap X$ . Put  $\mathcal{E} = Z(K) \cap X$  and  $\mathcal{D} = Z(K^*) \cap X$ . From ([7], 4.2) we have  $Q(X, K) = v(X, \mathcal{E})$  and  $Q(X, K^*) = v(X, \mathcal{D})$ , therefore  $v(X, \mathcal{E}) = v(X, \mathcal{D})$ . According to (1) and (2) it follows that  $\hat{\mathcal{D}} = \hat{\mathcal{E}}$  and since  $\mathcal{D} \subset \mathcal{E}$ ,  $\mathcal{D} \neq \mathcal{E}$ , the base  $\mathcal{D}$  is not complete. ■

**EXAMPLE 3.** Let  $Y$  be a locally compact, pseudocompact space which is not compact. Let  $X$  be the product space of  $Y$  with the two-point discrete space  $\{1, 2\}$ . If  $C_i = (\beta Y \sim Y) \times \{i\}$ ,  $1 \leq i \leq 2$ , then  $C_1$  and  $C_2$  are disjoint, homeomorphic, compact subsets of  $\beta X \sim X$ . However, since  $X$  is pseudocompact we have  $vX = \beta X$ , and consequently each base on  $X$  is complete.

**COROLLARY 4.** *A paracompact space  $X$  is Lindelöf if and only if each base on  $X$  is complete.*

**PROOF.** *Necessity.* By ([15], Lemma 2.7), if  $X$  is Lindelöf the only base on  $X$  is  $Z(X)$ . *Sufficiency.* Assume that  $X$  is not Lindelöf. Then, by a Theorem of Michael [10], there is an uncountable discrete family of open sets in  $X$ . From this family we can obtain disjoint, closed, discrete subsets  $C_1, C_2$  of  $X$  each of which has cardinal  $\aleph_1$ . Since  $X$  is normal  $C_1$  and  $C_2$  are completely separated and therefore  $\text{cl}_{\beta X} C_1 \cap \text{cl}_{\beta X} C_2 = \emptyset$ . If  $K_i = \text{cl}_{\beta X} C_i \sim C_i$ , then  $K_i$  is a compact subset of  $\beta X \sim X$  which is not  $z$ -separated from  $X$  because the cardinality of  $C_i$  is  $\aleph_1$ . On the other hand, as  $C_i$  is  $C$ -embedded in  $X$  we have  $\text{cl}_{\beta X} C_i = \beta C_i$  and  $\text{cl}_{vX} C_i = vC_i = C_i$ . Consequently  $K_i \cap vX = \emptyset$ ,  $1 \leq i \leq 2$ , and  $K_1$  is homeomorphic to  $K_2$ . According to Theorem 2 there exists a noncomplete base on  $X$ . ■

We need the following Lemma whose proof is left to the reader.



LEMMA 5. Suppose that  $X \subset \bar{X}$  and  $Y \subset \bar{Y}$ . The  $Q$ -closure of  $X \times Y$  in  $\bar{X} \times \bar{Y}$  is the product  $Q(X, \bar{X}) \times Q(Y, \bar{Y})$ .

PROPOSITION 6. Let  $Y$  be a space having at least two points. If each base on  $X \times Y$  is complete, then  $vX$  is Lindelöf.

PROOF. Assume that  $vX$  is not Lindelöf. Then there exists a compact set  $K$  in  $\beta X \sim vX$  which is not  $z$ -separated from  $X$ . Let  $p_1, p_2$  be distinct points of  $Y$  and let  $C_i = K \times \{p_i\}$ ,  $1 \leq i \leq 2$ . By Lemma 5  $vX \times vY = Q(X \times Y, \beta X \times \beta Y)$ , hence  $C_1$  and  $C_2$  are homeomorphic, disjoint, compact subsets of  $(\beta X \times \beta Y) \sim Q(X \times Y, \beta X \times \beta Y)$ .

We show next that  $C_2$  is not  $z$ -separated from  $X \times Y$ . Let  $Z$  be a zero-set in  $\beta X \times \beta Y$  containing  $C_2$ . Since each zero-set is a  $G_\delta$ , there exists a sequence  $\{U_n: n \in \mathbb{N}\}$  of open sets in  $\beta X \times \beta Y$  such that  $Z = \bigcap \{U_n: n \in \mathbb{N}\}$ . Each  $U_n$  is a neighborhood of  $C_2 = K \times \{p_2\}$ , hence there is a zero-set  $Z_n$  in  $\beta X$  such that  $K \subset Z_n$  and  $Z_n \times \{p_2\} \subset U_n$ . Then the set  $Z_0 = \bigcap \{Z_n: n \in \mathbb{N}\}$  is a zero-set in  $\beta X$  which contains  $K$ . As  $K$  is not  $z$ -separated from  $X$ , there exists a point  $x \in Z_0 \cap X$ . Now it is easy to check that  $(x, p_2) \in Z \cap (X \times Y)$ .

The proof is concluded by applying Theorem 2 to  $\beta X \times \beta Y$ . ■

THEOREM 7. For every space  $X$ , the following statements are equivalent:

- (1)  $vX$  is Lindelöf.
- (2) For every compact space  $K$  of non-measurable cardinal having at least two points, each base on  $X \times K$  is complete.
- (3) There exists a compact space  $K$  of non-measurable cardinal having at least two points such that each base on  $X \times K$  is complete.

PROOF. (1) implies (2). If  $K$  is a compact space of nonmeasurable cardinal then  $v(X \times K) = vX \times K$  ([4], Theorem 5.3). Since  $vX$  is Lindelöf, it follows that  $v(X \times K)$  is Lindelöf and therefore each base on  $X \times K$  is complete. (2) implies (3) Obvious. (3) implies (1) It is a consequence of Proposition 6. ■

The above result enables us to obtain easily examples of noncomplete bases on connected separable spaces.

EXAMPLE 8. Let  $X$  be the product space  $\mathbb{R}^c$ . Then  $X$  is separable, connected and realcompact. Because an uncountable product of copies of  $\mathbb{N}$  fails to be normal [16], it follows that  $X$  is not normal. According to Theorem 7, there exists a noncomplete base on the connected separable space  $X \times \beta X$ .

From Proposition 6 if each base on  $X \times X$  is complete then  $vX$  is Lindelöf. The converse is not true as the next example shows. We write  $\mathbb{R}$  for the set of all real numbers provided with the usual topology.

EXAMPLE 9. Let  $X$  be the set  $\mathbb{R}$  provided with the topology for which a base for the open sets is the family of all sets of the form  $[a, b)$ , where  $a, b \in \mathbb{R}$ . It is known that  $X$  is a Lindelöf space. Let  $Y$  be the product space  $X \times X$  and consider the sets  $S_1 = \{(x, y) \in Y: x + y = 1\}$ ,  $S_2 = \{(x, y) \in Y: x + y = -1\}$ . The map  $\sigma: Y \rightarrow Y$  defined  $\sigma(x, y) = (x - 1, y - 1)$  is an automorphism of  $Y$  which can be extended to an automorphism  $\phi$  of  $\beta Y$ . It is easy to see that the restriction of  $\phi$  to  $C_1 = \text{cl}_{\beta Y} S_1 \sim S_1$  is a homeomorphism from  $C_1$  onto  $C_2 = \text{cl}_{\beta Y} S_2 \sim S_2$ . Since  $Y$  is realcompact we have  $Y = vY = Q(Y, \beta Y)$ , therefore  $C_1$  and  $C_2$  are compact subsets of  $\beta Y \sim Q(Y, \beta Y)$ .



Moreover  $C_2$  is not  $z$ -separated from  $Y$  because  $S_2$  is an uncountable, discrete, closed subset of  $Y$ . On the other hand, since  $S_1$  and  $S_2$  are disjoint zero-sets in  $Y$ , it follows that  $\text{cl}_{\beta X} S_1 \cap \text{cl}_{\beta Y} S_2 = \emptyset$ . Hence  $C_1 \cap C_2 = \emptyset$ . By applying Theorem 2 to  $\beta Y$  we obtain a noncomplete base on  $Y$ .

### Inverse-closed subalgebras of $C(X)$

By an algebra on  $X$  is meant a subalgebra of  $C(X)$  which separates points and closed sets, contains the constants and is closed under inversion and uniform convergence. If  $A$  is an algebra on  $X$  and  $Z(A) = \{Z(f) : f \in A\}$ , the map  $A \rightarrow Z(A)$  is a one-to-one correspondence between the family of all algebras on  $X$  and the family of all bases on  $X$  ([15], Theorem 4.3). Moreover, if  $A$  is an algebra on  $X$  isomorphic to  $C(Y)$  for some space  $Y$ , then  $\nu Y = \nu(X, Z(A))$  and  $\beta Y = \omega(X, Z(A))$ . Therefore an algebra  $A$  on  $X$  is isomorphic to  $C(Y)$  for some space  $Y$  if and only if  $Z(A)$  is a complete base on  $X$  ([2], Theorem 5).

From Corollary 4, Theorem 7 and the above equivalence we deduce the following results

**THEOREM 10.** *A paracompact space  $X$  is Lindelöf if and only if each algebra on  $X$  is a  $C(Y)$ .*

**THEOREM 11.** *For every space  $X$ , the following statements are equivalent:*

- (1)  $\nu X$  is Lindelöf.
- (2) For every compact space  $K$  of nonmeasurable cardinal having at least two points, each algebra on  $X \times K$  is a  $C(Y)$ .
- (3) There is a compact space  $K$  of nonmeasurable cardinal having at least two points such that each algebra on  $X \times K$  is a  $C(Y)$ .

### REFERENCES

- [1] ALÖ, R. A. and SHAPIRO, H. L., *Normal Topological Spaces*, Cambridge Tracts in Mathematics, No. 65. Cambridge University Press, New York—London, 1974. MR 52 # 11808.
- [2] BLASCO, J. L., Complete bases and Wallman realcompactifications, *Proc. Amer. Math. Soc.* 75 (1979), 114—118. MR 80g: 54033.
- [3] BLASCO, J. L., Noncomplete bases on discrete spaces, *Acta Math. Acad. Sci. Hungar.* 36 (1980), 115—117. MR 82j: 54044.
- [4] COMFORT W. W. and NEGREPONTIS, S., Extending continuous functions on  $X \times Y$  to subsets of  $\beta X \times \beta Y$ , *Fund. Math.* 59 (1966), 1—12. MR 37 # 782.
- [5] GILLMAN, L. and JERISON, M., *Rings of continuous functions*, The University Series in Higher Mathematics, D. Van Nostrand Co., Princeton, N. J.—Toronto—London—New York, 1960. MR 22 # 6994.
- [6] HAGER, A. W. and JOHNSON, D. G., A note on certain subalgebras of  $C(X)$ , *Canad. J. Math.* 20 (1968), 389—393. MR 36 # 5697.
- [7] HAGER, A. W., On inverse-closed subalgebras of  $C(X)$ , *Proc. London Math. Soc.* (3) 19 (1969), 233—257. MR 39 # 6261.
- [8] HENRIKSEN, M. and JOHNSON, D. G., On the structure of a class of archimedean lattice-ordered algebras, *Fund. Math.* 50 (1961/62), 73—94. MR 24 # A3524.
- [9] ISBELL, J. R., Algebras of uniformly continuous functions, *Ann. of Math.* (2) 68 (1958), 96—125. MR 21 # 2177.
- [10] MICHAEL, E. A., A note on paracompact spaces, *Proc. Amer. Math. Soc.* 4 (1953), 831—833. MR 15—144.

- [11] MRÓWKA, S., Some set-theoretic constructions in topology, *Fund. Math.* **94** (1977), 83—92. *MR* **55** # 6364.
- [12] SMIRNOV, YU M., On normally disposed sets of normal spaces, *Mat. Sbornik* N. S. **29 (71)** (1951), 173—176. (in Russian). *MR* **13**—371.
- [13] STEINER, E. F., Normal families and completely regular spaces, *Duke Math. J.* **33** (1966), 743—745. *MR* **33** # 7975.
- [14] STEINER, E. F., Wallman spaces and compactifications, *Fund. Math.* **61** (1967/68), 295—304. *MR* **36** # 5899.
- [15] STEINER, A. K. and STEINER, E. F., Nest generated intersection rings in Tychonoff spaces, *Trans. Amer. Math. Soc.* **148** (1970), 589—601. *MR* **41** # 7637.
- [16] STONE, A. H., Paracompactness and product spaces, *Bull. Amer. Math. Soc.* **54** (1948), 977—982. *MR* **10**—204.

(Received September 14, 1983)

CÁTEDRA DE MATEMÁTICAS II  
FACULTAD DE CIENCIAS MATEMÁTICAS  
UNIVERSIDAD DE VALENCIA  
DOCTOR MOLINER, 50  
BURJASOT  
VALENCIA  
SPAIN

# NUMBER BASES IN QUADRATIC FIELDS

EDWARD H. GROSSMAN

Let  $K$  be a number field,  $O_K$  the ring of integers and suppose  $O_K = \mathbb{Z}[\theta]$ . For  $\alpha \in K$  we let  $N(\alpha)$  and  $Tr(\alpha)$  denote the norm and trace from  $K$  to  $\mathbb{Q}$ . If every  $\alpha \in O_K$  can be represented in the form

$$(1) \quad \alpha = r_0 + r_1\theta + \dots + r_l\theta^l$$

where each  $r_i$  satisfies  $0 \leq r_i < |N(\theta)|$  we call  $\theta$  a base. For the rational field  $K = \mathbb{Q}$  the only bases are  $\theta = -A$ ,  $A \geq 2$ . Kátai and Szabó ([3]) and Kátai and Kovács ([1] and [2]) determined all bases in quadratic extensions. For convenience we have reformulated the latter's results as follows:

**THEOREM 1.** *Let  $[K: \mathbb{Q}] = 2$ . An element  $\theta \in O_K$  satisfying  $O_K = \mathbb{Z}[\theta]$  is a base if and only if the following conditions hold:*

$$(i) \quad Tr(\theta) \leq 1 \quad (ii) \quad |Tr(\theta)| \leq N(\theta) \quad (iii) \quad 2 \leq N(\theta).$$

Let  $\theta$  be a base and denote by  $l(\alpha, \theta)$  the exponent  $l$  in (1). The methods used to prove Theorem 1 yield very weak estimates for  $l$ . In this paper a different approach is used which yields the following considerably sharper results.

**THEOREM 2.** *Let  $K$  be an imaginary quadratic field and  $\theta$  a base. Then*

$$l(\alpha, \theta) = \left\lfloor \frac{\log N(\alpha)}{\log N(\theta)} \right\rfloor + O(1)$$

where  $O(1)$  is a constant depending only on  $K$ .

For a real quadratic field  $K = \mathbb{Q}(\sqrt{D})$ ,  $D > 0$ , let  $M(\alpha) = (\alpha^2 + \bar{\alpha}^2)/2$ , where  $\bar{\alpha}$  is the conjugate of  $\alpha$  and let  $m(\alpha) = \min(\alpha^2, \bar{\alpha}^2)$ . We prove

**THEOREM 3.** *Let  $K$  be a real quadratic field and  $\theta$  a base. Then*

$$l(\alpha, \theta) \leq \left\lfloor \frac{\log M(\alpha)}{\log \frac{m(\theta)}{1 + (D+2)^{-2}}} \right\rfloor + O(1)$$

where  $O(1)$  depends only on  $K$ .

To facilitate the proofs of the theorems we introduce the notion of a congruence chain mod  $\theta$ .

**DEFINITION 1.** A finite sequence  $\alpha_0, \alpha_1, \dots, \alpha_n$  is a congruence chain mod  $\theta$  of length  $n$  if for every  $i$ ,  $1 \leq i \leq n$ ,  $\alpha_i = (\alpha_{i-1} - r_{i-1})/\theta$ , where  $\alpha_{i-1} \equiv r_{i-1} \pmod{\theta}$  and  $r_{i-1} \in \mathbb{Z}$  satisfies  $0 \leq r_{i-1} < |N(\theta)|$ . Note that if  $O_K = \mathbb{Z}[\theta]$  then every  $\alpha \in O_K$  gives rise to a congruence chain mod  $\theta$ . Clearly (1) is equivalent to the statement that every sufficiently long congruence chain yields  $\alpha_n = 0$  and Theorems 2 and 3 provide estimates of the maximum length of a congruence chain with  $\alpha_n \neq 0$ .

**LEMMA 1.** Let  $K$  be a quadratic extension,  $\theta$  any element of  $O_K$  satisfying the conditions (i), (ii), (iii) of Theorem 1. If  $\bar{\theta}$  denotes the conjugate of  $\theta$  and  $\alpha = -\bar{\theta}$  or  $\alpha = 1 - \bar{\theta}$  then  $\alpha$  can be represented in the form (1) with  $l \leq 4$ .

**PROOF.** The verification of the following list is omitted. ( $\text{Tr}\theta$  and  $N\theta$  stand for  $\text{Tr}(\theta)$  and  $N(\theta)$  resp.)

$$\begin{aligned}
 -\bar{\theta} &= \begin{cases} -\text{Tr}\theta + \theta & \text{if } \text{Tr}\theta \equiv 0, \quad -\text{Tr}\theta < N\theta \\ \theta + (-1 - \text{Tr}\theta)\theta^2 + \theta^3 & \text{if } \text{Tr}\theta \equiv 0, \quad -\text{Tr}\theta = N\theta \\ (N\theta - \text{Tr}\theta) + \theta^2 & \text{if } \text{Tr}\theta = 1 \end{cases} \\
 1 - \bar{\theta} &= \begin{cases} (1 - \text{Tr}\theta) + \theta & \text{if } \text{Tr}\theta \equiv 0, \quad 1 - \text{Tr}\theta < N\theta \\ 2\theta + (-1 - \text{Tr}\theta)\theta^2 + \theta^3 & \text{if } \text{Tr}\theta \equiv 0, \quad 1 - \text{Tr}\theta = N\theta \equiv 3 \\ \theta^2 + \theta^3 + \theta^4 & \text{if } \text{Tr}\theta \equiv 0, \quad 1 - \text{Tr}\theta = N\theta = 2 \\ 1 + \theta + (-1 - \text{Tr}\theta)\theta^2 + \theta^3 & \text{if } \text{Tr}\theta \equiv 0, \quad -\text{Tr}\theta = N\theta \\ \theta & \text{if } \text{Tr}\theta = 1. \end{cases}
 \end{aligned}$$

We let  $\omega = \sqrt{D}$  if  $D \equiv 2, 3 \pmod{4}$  or  $\omega = (1 + \sqrt{D})/2$  if  $D \equiv 1 \pmod{4}$  be the canonical generator of  $O_K$  for  $K = \mathbb{Q}(\sqrt{D})$ . The roots are taken to be positive if  $D > 0$ .

**LEMMA 2.** Let  $\theta = -A + \omega$  satisfy  $N(\theta) > 0$ . If  $\alpha_0, \dots, \alpha_n$ ,  $n \geq 2$  is a congruence chain mod  $\theta$  where each  $\alpha_k = a_k + b_k \omega$  satisfies

$$(2) \quad |a_k + b_k A| < cN(\theta)$$

for some integer constant  $c$ , then for  $k \geq 2$

$$(3) \quad \alpha_k = a_k^* - b_k^* \bar{\theta} \quad \text{where} \quad 1 - c \leq a_k^* \leq c, \quad 1 - c \leq b_k^* \leq c.$$

**PROOF.** Let  $\alpha_{i-1}, \alpha_i, \alpha_{i+1}$  be any three consecutive terms in the congruence chain. Then from (2) it follows that  $r_{i-1} = a_{i-1} + b_{i-1}A + aN(\theta)$  where  $1 - c \leq a \leq c$  so  $\alpha_i = (\alpha_{i-1} - r_{i-1})/\theta = b_{i-1} - a\bar{\theta} = a_i + a\omega$ . Repeating this argument gives  $\alpha_{i+1} = a - b\bar{\theta}$  where  $a$  and  $b$  satisfy the conditions in (3).

Our next lemma shows that in an imaginary  $K$  the terms of a congruence chain eventually become "small".

**LEMMA 3.** Let  $K$  be imaginary,  $N(\theta) \geq 2$  and suppose  $\alpha = \alpha_0, \dots, \alpha_n$  is a congruence chain mod  $\theta$  with  $n \geq n_0 = [\log N(\alpha)/\log N(\theta)] + 11$ . Then

$$(4) \quad N(\alpha_n) < 3N(\theta).$$

PROOF. For every  $i \leq n$  we have

$$|\alpha_i| = \frac{|\alpha_{i-1} - r_{i-1}|}{|\theta|} \leq \frac{|\alpha_{i-1}|}{|\theta|} + \frac{N(\theta) - 1}{|\theta|}.$$

Iterating this inequality gives

$$|\alpha_n| \leq \frac{|\alpha_0|}{|\theta|^n} + (N(\theta) - 1) \sum_{i=1}^n |\theta|^{-i} < \frac{|\alpha_0|}{|\theta|^n} + |\theta| + 1 \leq |\theta| + \sqrt{2}(\sqrt{3} - 1)$$

if  $n \geq n_0$ , which, since  $|\theta| \geq \sqrt{2}$ , implies (4).

We need an appropriate generalization of Lemma 3 for real fields. Recall that for real  $\alpha$ ,  $M(\alpha) = (\alpha^2 + \bar{\alpha}^2)/2$  and  $m(\alpha) = \min(\alpha^2, \bar{\alpha}^2)$ . In what follows  $c_1, c_2, \dots, c_7$  will denote positive integer constants depending only on a particular field  $K$ .

Note first that the property of being a base as well as conditions (i), (ii) and (iii) are valid for  $\theta$  if and only if they hold for the conjugate  $\bar{\theta}$  of  $\theta$ . Moreover since  $\mathbf{Z}[\theta] = O_K$ ,  $\theta = -A \pm \omega$  and we may assume, using  $\bar{\theta}$  if necessary, that in fact  $\theta = -A + \omega$ .

LEMMA 4. Let  $\theta = -A + \omega$  be a base for the real field  $K = \mathbf{Q}(\sqrt{D})$ ,  $D > 0$  and suppose  $\alpha = \alpha_0, \dots, \alpha_n$  is a congruence chain mod  $\theta$ . If  $\Delta = (1 + (D + 2)^{-2})/m(\theta)$  then for  $n > n_0 = [\log M(\alpha)/\log(\Delta^{-1})] + 1$

$$(5) \quad M(\alpha_n) \leq c_1(N(\theta) - 1)^2/m(\theta).$$

PROOF. Note that  $M(a + b\sqrt{D}) = a^2 + b^2D$  for any rational  $a$  and  $b$ , from which it follows that if  $\alpha = a_1 + b_1\sqrt{D}$ ,  $\beta = a_2 + b_2\sqrt{D}$  are any elements of  $K$

$$M(\alpha\beta) = M(\alpha)M(\beta) + 4Da_1b_1a_2b_2.$$

Since  $2a_1b_1\sqrt{D} = M(\alpha) - \bar{\alpha}^2$  this yields that

$$(6) \quad M(\alpha\beta) \leq M(\alpha)(M(\beta) + 2|a_2b_2|\sqrt{D}).$$

Applying (6) to the terms of a congruence chain we obtain

$$M(\alpha_i) = M((\alpha_{i-1} - r_{i-1})\theta^{-1}) \leq M(\alpha_{i-1} - r_{i-1})/m(\theta).$$

Using that

$$M(\alpha_{i-1} - r_{i-1}) \leq (1 + (D + 2)^{-2})M(\alpha_{i-1}) + (N(\theta) - 1)^2c_2,$$

where  $c_2 = 1 + (D + 2)^2$ , gives

$$M(\alpha_i) \leq \Delta M(\alpha_{i-1}) + c_2(N(\theta) - 1)^2/m(\theta).$$

Iterating this we have

$$(7) \quad M(\alpha_n) \leq \Delta^n M(\alpha_0) + c_2(N(\theta) - 1)^2 \sum_{i=1}^n \Delta^{i-1}/m(\theta).$$

By condition (i) and (ii) of Theorem 1,  $\theta = -A + \omega$  satisfies  $A \geq 1 + \sqrt{D + 1}$  if  $D \equiv 2, 3 \pmod{4}$  or  $A \geq (3 + \sqrt{D + 4})/2$  if  $D \equiv 1 \pmod{4}$ . In either case

$$m(\theta) = (A - \omega)^2 > (1 + (D + 2)^{-1})^2 > 1 + (D + 2)^{-2} > 1$$

and  $\Delta = (1 + (D+2)^{-2})/m(\theta) < 1$ . From this (5) follows from (7) with  $c_1 = (c_2 + 1)(D+2)$ .

LEMMA 5. Let  $K = \mathbf{Q}(\sqrt{D})$  be any quadratic extension. If  $\theta = -A + \omega$  is a base and  $\alpha_0, \dots, \alpha_l$  is a congruence chain mod  $\theta$  where each  $\alpha_k = a_k + b_k\omega$  satisfies (2) then there is an integer  $l_0(c, K)$  such that  $\alpha_l = 0$  if  $l > l_0$ .

PROOF. Assume  $l > 4$ . Lemma 2 implies that for  $k \geq 2$

$$(8) \quad \alpha_k = a_k^* - b_k^* \bar{\theta} = (a_k^* + b_k^* A - b_k^* \text{Tr} \omega) + b_k^* \omega = a_k + b_k \omega$$

where  $|a_k^*| \leq c$ ,  $|b_k^*| \leq c$  and therefore for  $k \geq 2$ ,  $|a_k + b_k A| < c_3 |A|$ . If  $|A| > c_4$  this implies that for  $k \geq 2$ ,  $|a_k + b_k A| < N(\theta)$  so that by Lemmas 2 and 1,  $\alpha_k = 0$  for  $k \geq 9$ .

For a given field  $K$  there are only finitely many number bases with  $|A| \leq c_4$ . For each of these (8) yields a finite set of values for  $\alpha_2$ . If we let  $\bar{l} = \max$  of  $l(\alpha_2, \theta)$  for these  $\alpha_2$  then the lemma is proven with  $l_0 = \max(9, \bar{l})$ .

PROOF of Theorems 2 and 3. Let  $n_0$  be defined as in Lemma 3 and 4 for imaginary and real fields respectively. Suppose we have a congruence chain  $\alpha = \alpha_0, \dots, \alpha_l$  and  $l > n_0$ .

If  $K = \mathbf{Q}(\sqrt{D})$  is imaginary then by Lemma 3 for  $l > n_0$ ,  $N(\alpha_l) < 3N(\theta)$ . Writing  $\alpha_l = a_l + b_l \omega$  and  $\theta = -A + \omega$ , Cauchy's inequality gives  $|a_l + b_l A| < 4N(\theta)$ . Hence Lemma 5 implies that  $\alpha_l = 0$  if  $l > n_0 + l_0$ , where  $l_0$  depends only on  $K$ . This establishes the upper bound implied in Theorem 2, the lower bound being clear.

For real fields, Lemma 4 gives for  $l > n_0$

$$(9) \quad |a_l + b_l A| \leq \frac{2(M(\alpha)M(\theta))^{1/2}}{\sqrt{D}} < c_5 \frac{M(\theta)^{1/2}}{m(\theta)^{1/2}} N(\theta).$$

If  $A > c_6$  then (9) implies  $|a_l + b_l A| < c_6 N(\theta)$  which again by Lemma 5 yields  $\alpha_l = 0$  for  $l > n_0 + l_0$ . If  $A \leq c_6$  then  $\theta$  can have only finitely many values and for each of these we have by Lemma 4 that for  $l > n_0$ ,  $M(\alpha_l) \leq c_1(N(\theta) - 1)^2/m(\theta) < c_7$ . Thus the  $\alpha_{n_0+1}$  belong to a finite set and so there is a maximum exponent  $l_0^*$  in the representation (1) for these elements. Thus in the real case a representation (1) exists with  $l \leq n_0 + \max(l_0, l_0^*)$  and the proof of Theorem 3 is complete.

#### REFERENCES

- [1] KÁTAI, I. and KOVÁCS, B., Kanonische Zahlensysteme in der Theorie der quadratischen algebraischen Zahlen, *Acta Sci. Math. (Szeged)* **42** (1980), 99—107. MR **81i**: 12002.
- [2] KÁTAI, I. and KOVÁCS, B., Canonical number systems in imaginary quadratic fields, *Acta Math. Acad. Sci. Hungar.* **37** (1981), 159—164. MR **83a**: 12005.
- [3] KÁTAI, I. and SZABÓ, J., Canonical number systems for complex integers, *Acta Sci. Math. (Szeged)* **37** (1975), 255—260. MR **52** # 10590.

(Received October 5, 1983)

DEPARTMENT OF MATHEMATICS  
CITY COLLEGE OF NEW YORK  
CONVENT AVENUE AT 138TH STREET  
NEW YORK, NY 10031  
U.S.A.



# ON RINGS WITH MODIFIED CHAIN CONDITIONS

DINH VAN HUYNH

1. Throughout the present paper we consider associative rings  $A$  with Jacobson radical  $J(A)$ . In [1] and [2] the *almost artinian rings* and *modules* have been introduced and investigated. A ring  $A$  is said to be almost right artinian, if for every infinite descending chain  $R_1 \supseteq R_2 \supseteq \dots$  of right ideals  $R_i$  of  $A$  there are positive integers  $m, q$  such that  $R_m A^q \subseteq R_i$  holds for all  $i$ , or equivalently there exists a positive integer  $p$  such that  $R_p A^p \subseteq R_i$  holds for all  $i$ . Thus every nilpotent ring is almost right and left artinian. As main result of this paper we will prove the following

THEOREM 1. *For a ring  $A$  the following are equivalent:*

- (i)  *$A$  is almost right artinian,*
- (ii)  *$A$  contains a right artinian left ideal  $B$  with a right identity  $e$  and a nilpotent left ideal  $C$  with  $Ce = (0)$  such that  $A$  is a groupdirect sum of  $B$  and  $C$ :  $A = B \oplus C$ .*

COROLLARY 2. (1) *For every almost right artinian ring  $A$ ,  $A^m$  is also almost right artinian ( $m=1, 2, \dots$ ).*

(2) *Every ideal of a (almost) right artinian ring is an almost right artinian ring.*

In section 3 we consider almost MHR-rings and prove that for every almost MHR-ring  $A$ ,  $J(A)$  is a nil ideal and  $A/J(A)$  is a MHR-ring.

2. PROOF of Theorem 1. (i) $\Rightarrow$ (ii). Let  $A$  be an almost right artinian ring. Then there exists an idempotent  $e$  in  $A$  such that  $A(1-e) \stackrel{\text{def}}{=} \{a - ae, a \in A\}$  is contained in  $J(A)$  (cf. [1] or [2]). Noting

$$(1) \quad A = Ae \oplus A(1-e),$$

and considering a descending chain of right ideals  $R_i$  ( $i=1, 2, \dots$ ) of  $Ae$ :

$$(2) \quad R_1 \supseteq R_2 \supseteq \dots,$$

we have that

$$R_1 + R_1 A \supseteq R_2 + R_2 A \supseteq \dots$$

is a descending chain of right ideals of  $A$ . By assumption there is a positive integer  $p$  with

$$(3) \quad (R_p + R_p A) A^p = R_p A^p + R_p A^{p+1} \subseteq R_i + R_i A$$

for all  $i$ . Using (1) and by multiplying (3) by  $e$  on the right we get  $R_p Ae + R_p Ae \subseteq$

$\subseteq R_i Ae + R_i = R_i$ , therefore  $R_p$  is contained in all  $R_i$ . By (2),  $R_p = R_{p+1} = R_{p+2} = \dots$ , proving that  $B \stackrel{\text{def}}{=} Ae$  is a right artinian ring. Since  $J(A)$  is nilpotent (cf. [1], [2]),  $C \stackrel{\text{def}}{=} A(1-e)$  is nilpotent.

(ii)  $\Rightarrow$  (i). Let  $n$  be a positive integer with  $C^n = (0)$ ,  $C^{n-1} \neq (0)$ , and

$$(4) \quad S_1 \supseteq S_2 \supseteq \dots$$

be an infinite descending chain of right ideals  $S_i$  of  $A$ . Since  $Ae = B$  ( $e$  is a right identity of  $B$ ) is a right artinian ring, there is a positive integer  $m$  with

$$(5) \quad S_m e = S_{m+1} e = \dots$$

Noting  $S_i = S_i e \oplus S_i(1-e)$ , it follows that  $S_i(1-e)$  is contained in  $C$ . Clearly  $A^p = B \oplus BC$  holds for all  $p \geq n$ . Hence

$$S_m A^p = (S_m e \oplus S_m(1-e))(B \oplus BC) \subseteq S_m e \oplus S_m BC.$$

By (5),  $S_m e \oplus S_m BC = S_{m+j} e \oplus S_{m+j} e \cdot C \subseteq S_{m+j}$  for all  $j = 1, 2, \dots$ . Hence  $S_m A^p$  is contained in all  $S_i$ , proving (i).

The Theorem is proved.

PROOF of Corollary 2. For the case that  $A$  is an almost right artinian ring, (1) in fact is a special case of (2). Hence it is enough to prove (2). Let  $A$  be an almost right artinian ring, and  $I$  be an ideal of  $A$ . Then by Theorem 1,  $I/J(I)$  has a right identity  $\bar{e}$ , because  $J(I) = I \cap J(A)$ . Let  $e$  be an idempotent of  $I$  with  $e \in \bar{e}$ . Then

$$(6) \quad I = Ie \oplus I(1-e) = Ae \oplus I(1-e),$$

where  $I(1-e)$  is a nilpotent left ideal of  $I$ . By the same way as in the proof of Theorem 1 (the part (i)  $\Rightarrow$  (ii)) we can verify that  $Ae$  ( $= Ie$ ) is a right artinian ring. Now, (6) and Theorem 1 show that  $I$  is an almost right artinian ring.

By Corollary 2 one can use the results obtained about almost right artinian rings for studying the structure of ideals of right artinian rings.

3. Following [3], [4], a ring  $A$  is called a MHR-ring, if  $A$  satisfies minimal condition on principal right ideals. Now, a ring  $A$  is defined to be almost MHR-ring, if for every infinite descending chain  $R_1 \supseteq R_2 \supseteq \dots$  of principal right ideals  $R_i$  of  $A$  there exist positive integers  $m, q$  such that  $R_m A^q \subseteq R_i$  holds for all  $i$ , or equivalently, there exists a positive integer  $p$  such that  $R_p A^p \subseteq R_i$  holds for all  $i$ .

THEOREM 3. For every almost MHR-ring  $A$ ,  $J(A)$  is a nil ring and the factor ring  $A/J(A)$  is a MHR-ring.

PROOF. Let  $A$  be an almost MHR-ring, and  $x$  is an element of  $J(A)$ . Denote by  $(x^i)_r$  the principal right ideal of  $A$  generated by  $x^i$  ( $i = 1, 2, \dots$ ). Then  $(x)_r \supseteq (x^2)_r \supseteq \dots$ . By assumption there is a positive integer  $m$  such that  $(x^m)_r A^m$  is contained in all  $(x^i)_r$ . Since  $(x^m)_r A^m = x^m A^m + x^m A \cdot A^m = x^m A^m$ ,  $y \stackrel{\text{def}}{=} x^{2m}$  is contained in  $(x^m)_r A^m$ , therefore  $x^m A^m = (x^{2m})_r = (x^{2m+1})_r = \dots$ . Then there are an integer  $h$  and an element  $a$  of  $A$  with  $y = yxh + y(xa) = y(xh + xa) = ys$ , where  $s \stackrel{\text{def}}{=} xh + xa \in J(A)$ . As is

well-known, for  $s$  as an element of  $J(A)$  there is a  $t \in J(A)$  with  $s - st + t = 0$ . Hence

$$x^{2m} = y = y - y(s - st + t) = (y - ys) - (y - ys)t = 0,$$

proving that  $J(A)$  is a nil ideal of  $A$ .

It is obvious, that  $A/J(A)$  is an almost MHR-ring. Hence, for proving that  $A/J(A)$  is a MHR-ring, it is enough to consider the case  $J(A) = (0)$ . Firstly we prove that every non-zero right ideal  $R$  of  $A$  contains an idempotent minimal right ideal of  $A$ . Let  $0 \neq x_1 \in R$ . Since  $J(A) = (0)$ ,  $SA \neq (0)$  holds for each non-zero right ideal  $S$  of  $A$ . If  $(x_1)_r A$  is not minimal, we can find an  $0 \neq x_2 \in (x_1)_r A$  with  $(x_1)_r A \supset (x_2)_r$ . If  $(x_2)_r A^2$  is not minimal, we can find an  $0 \neq x_3 \in (x_2)_r A^2$  such that  $(x_2)_r A^2 \supset (x_3)_r$ . Successively we find a descending chain of principal right ideals of  $A$

$$(1) \quad (x_1)_r \supset (x_2)_r \supset (x_3)_r \supset \dots$$

with  $(x_i)_r A^i \not\subseteq (x_{i+1})_r$ . By assumption the chain (1) must be finite, i.e. there exists an  $(x_m)_r$  in (1) such that  $(x_m)_r$  is minimal. Clearly  $(x_m)_r A = (x_m)_r$ , it can be generated by an idempotent, say  $e_1$ :  $e_1 A = (x_m)_r$ . Hence  $(x_1)_r = e_1 A \oplus (1 - e_1)(x_1)_r$ . If  $R_2 \stackrel{\text{def}}{=} (1 - e_1)(x_1)_r$  is nonzero, it contains an idempotent minimal right ideal  $e_2 A$  ( $e_2^2 = e_2 \in A$ ):  $(x_1)_r = e_1 A \oplus e_2 A \oplus (1 - e_2)R_2$ . By the same way we get  $(x_1)_r = e_1 A \oplus \dots \oplus e_i A \oplus (1 - e_i)R_i$ . One can easily verify that each  $(1 - e_i)R_i$  for each  $i$  is a principal right ideal with  $(1 - e_j)R_j A^j \not\subseteq (1 - e_{j+1})R_{j+1}$  if  $R_{j+1} \neq (0)$ . Since  $A$  is an almost MHR-ring, there must exist a positive integer  $m$  such that  $(1 - e_m)R_m$  is minimal. Hence  $(x_1)_r = e_1 A \oplus \dots \oplus e_m A$  is a direct sum of minimal right ideals of  $A$ . Note

$$A = \sum_{x \in A} (x)_r = \sum_i^{\oplus} e_i A$$

where each  $e_i A$  is an idempotent minimal right ideal of  $A$ , proving that  $A$  is a MHR-ring.

Theorem 3 is proved.

REMARK. We do not know, whether every ideal of a MHR-ring or every nil ring is an almost MHR-ring or not.

#### REFERENCES

- [1] CATER, F. S., Modified chain conditions for rings without identity, *Yokohama Math. J.* **27** (1979), 1—22. MR **81c**: 16019.
- [2] KOMATSU, H. and TOMINAGA, H., On modified chain conditions, *Math. J. Okayama Univ.* **22** (1980), 131—139. MR **82e**: 16008a.
- [3] SZÁSZ, F., Über Ringe mit Minimalbedingung für Hauptideale I, *Publ. Math. Debrecen*, **7** (1960), 54—64. MR **24** # A145.
- [4] SZÁSZ, F., Über Ringe mit Minimalbedingung für Hauptideale II, *Acta Math. Acad. Sci. Hungar.* **12** (1961), 417—439. MR **26** # 6207.

(Received November 10, 1983)



# ON THE LAW OF THE ITERATED LOGARITHM FOR RANDOMLY INDEXED PARTIAL SUMS WITH TWO APPLICATIONS

ALLAN GUT

An Anscombe-type of law of the iterated logarithm due to Huggins is slightly modified and applied to generalized (non-linear) renewal theory and to certain stopped sums.

## 1. Introduction

The starting point of this note is the following law of the iterated logarithm for randomly indexed partial sums.

**THEOREM 1.** *Let  $\{X_n\}_{n=1}^\infty$  be i.i.d. random variables with  $EX_1=0$  and  $EX_1^2=\sigma^2<\infty$  and set  $S_n=\sum_{k=1}^n X_k$ ,  $n=1, 2, \dots$ . Let  $\{b_n\}_{n=1}^\infty$  be a strictly increasing sequence of positive reals, increasing to infinity, such that*

$$(1.1) \quad \frac{b_{n+1}}{b_n} \rightarrow B, \quad 1 \leq B < \infty,$$

*and let  $\{\tau_n\}_{n=1}^\infty$  be a strictly increasing sequence of positive, integer valued random variables with  $\tau_1 \geq 3$  and such that*

$$(1.2) \quad \frac{\tau_n}{b_n} \xrightarrow{\text{a.s.}} \theta \quad \text{as } n \rightarrow \infty,$$

*where  $\theta$  is a positive, finite constant. Then, the set of cluster points of the sequence*

$$\left\{ \frac{S_{\tau_n}}{\sqrt{2\sigma^2 \tau_n \log \log \tau_n}}; n \geq 1 \right\}$$

*coincides with  $[-1, 1]$  a.s.*

The theorem, which can be considered as an “Anscombe-theorem” for the law of the iterated logarithm (cf. [1]), is contained in Huggins [13], Chow et al. [5] and Chang and Hsiung [2], in which more general versions of the law of the iterated logarithm are proved.

---

1980 *Mathematics Subject Classifications*. Primary 60F15, 60G40, 60G50; Secondary 60K05, 60K10.

*Key words and phrases*. i.i.d. random variables, randomly indexed partial sums, law of the iterated logarithm, regular variation, first passage times, stopped sums.

The purpose of this note is to extend Theorem 1 to continuous time and, in particular, to the case when the normalizing function (corresponding to  $\{b_n\}$  above) is regularly varying with positive exponent and, further, to apply this extension to two applications. In the first one, a different proof of a result of Chow and Hsiung [3], dealing with first passage times, will be obtained in a slightly more general setting and in the second one a law of the iterated logarithm for certain stopped sums will be given.

## 2. A continuous time extension

We use the notation  $C(\{x_n\})$  to denote the set of cluster points (the cluster set) in  $\mathbf{R}$  of a sequence  $\{x_n\}$ . Further, set  $\log_2 x = \max\{1, \log \log x\}$ ,  $x > 0$ .

As a first step we state a lemma, which will be used repeatedly without explicit reference and whose proof is immediate.

LEMMA 1. Let  $E \subset \mathbf{R}$  and suppose that  $\{Y(t); t > 0\}$  is a family of random variables such that  $C(\{Y(t); t > 0\}) = E$  a.s., Further, let  $\{\xi(t); t > 0\}$  and  $\{\eta(t); t > 0\}$  be families of random variables such that  $\xi(t) \xrightarrow{\text{a.s.}} 1$  and  $\eta(t) \xrightarrow{\text{a.s.}} 0$  as  $t \rightarrow \infty$ . Then

$$C(\{\xi(t)Y(t) + \eta(t); t > 0\}) = E \quad \text{a.s.}$$

As an immediate consequence of Theorem 1 and Lemma 1 we obtain

THEOREM 2. Under the assumptions of Theorem 1 we have

$$C\left(\left\{\frac{S_{\tau_n}}{\sqrt{2\sigma^2\theta b_n \log_2 b_n}}; n \geq 1\right\}\right) = [-1, 1] \quad \text{a.s.}$$

Next we shall extend this result to the case where  $\tau$  and  $b$  are functions of a continuous variable.

THEOREM 3. Let  $\{X_n\}_{n=1}^\infty$  be as before and suppose that  $b(t)$ ,  $t > 0$ , is a positive function, strictly increasing to infinity, such that

$$(2.1) \quad \frac{b(n+1)}{b(n)} \rightarrow B \quad \text{as } n \rightarrow \infty, \quad \text{where } 1 \leq B < \infty.$$

Further, let  $\{\tau(t); t > 0\}$  be a family of positive, integer valued random variables, strictly increasing in  $t$  and such that

$$(2.2) \quad \frac{\tau(t)}{b(t)} \xrightarrow{\text{a.s.}} 0 \quad \text{as } t \rightarrow \infty, \quad \text{where } 0 < \theta < \infty.$$

Then

$$(2.3) \quad C\left(\left\{\frac{S_{\tau(t)}}{\sqrt{2\sigma^2\theta b(t) \log_2 b(t)}}; t > 0\right\}\right) = [-1, 1] \quad \text{a.s.}$$

PROOF. From Theorem 2 we know that

$$C\left(\left\{\frac{S_{\tau(n)}}{\sqrt{2\sigma^2\theta b(n) \log_2 b(n)}}; n \geq 1\right\}\right) = [-1, 1] \quad \text{a.s.}$$



Furthermore, with probability 1, we have

$$\begin{aligned} C\left(\left\{\frac{S_{\tau(n)}}{\sqrt{2\sigma^2\theta b(n)\log_2 b(n)}}; n \geq 1\right\}\right) &\subset C\left(\left\{\frac{S_{\tau(t)}}{\sqrt{2\sigma^2\theta b(t)\log_2 b(t)}}; t > 0\right\}\right) = \\ &= C\left(\left\{\frac{S_{\tau(t)}}{\sqrt{2\sigma^2\tau(t)\log_2 \tau(t)}}; t > 0\right\}\right) \subset C\left(\left\{\frac{S_n}{\sqrt{2\sigma^2 n \log_2 n}}; n \geq 1\right\}\right) = [-1, 1] \quad \text{a.s.}, \end{aligned}$$

which proves the theorem.

### 3. Regularly varying normalizations

In this section we shall apply Theorem 3 to regularly varying normalizing functions  $b(t)$ ,  $t > 0$ . For some general facts about such functions we refer to Feller [6], Chapter VIII, de Haan [9] and Seneta [14].

**THEOREM 4.** Let  $\{X_n\}_{n=1}^\infty$  be i.i.d. random variables with  $EX_1=0$  and  $EX_1^2=\sigma^2<\infty$  and suppose that  $b(t)$ ,  $t>0$  is a positive, strictly increasing function which is regularly varying at infinity with exponent  $\varrho>0$ . Further, let  $\{\tau(t); t>0\}$  be a family of positive, integer valued random variables, strictly increasing in  $t$  and such that

$$(3.1) \quad \frac{\tau(t)}{b(t)} \xrightarrow{\text{a.s.}} \theta \quad \text{as } t \rightarrow \infty, \quad \text{where } 0 < \theta < \infty.$$

Then

$$(3.2) \quad C\left(\left\{\frac{S_{\tau(t)}}{\sqrt{2\sigma^2\theta b(t)\log_2 b(t)}}; t > 0\right\}\right) = [-1, 1] \quad \text{a.s.}$$

In particular,

$$(3.3) \quad \limsup_{t \rightarrow \infty} \frac{S_{\tau(t)}}{\sqrt{2\sigma^2\theta b(t)\log_2 b(t)}} = 1 \quad \text{a.s.}$$

Moreover,

$$(3.4) \quad C\left(\left\{\frac{S_{\tau(t)}}{\sqrt{2\sigma^2\theta b(t)\log_2 t}}; t > 0\right\}\right) = [-1, 1] \quad \text{a.s.},$$

in particular,

$$(3.5) \quad \limsup_{t \rightarrow \infty} \frac{S_{\tau(t)}}{\sqrt{2\sigma^2\theta b(t)\log_2 t}} = 1 \quad \text{a.s.}$$

**REMARK 3.1.** It is clear that the theorem remains true if the assumptions on  $b(t)$  only hold in some interval  $[A, \infty)$ .

**REMARK 3.2.** From the theory of regularly varying functions it follows that the "typical" examples one should have in mind are  $b(t)=t^\varrho$  and  $b(t)=t^\varrho \log t$ .

**PROOF.** (3.2) follows from Theorem 3 once we have shown that (2.1) is satisfied for some  $B$ ,  $1 \leq B < \infty$ .

Let  $\delta > 0$ . By monotonicity it follows that

$$b(t) \leq b(t+1) = b(t(1+t^{-1})) \leq b(t(1+\delta)) \quad \text{for } t > 1/\delta,$$

and so

$$1 \cong \liminf_{t \rightarrow \infty} b(t+1)/b(t) \cong \limsup_{t \rightarrow \infty} b(t+1)/b(t) \cong (1+\delta)^e,$$

the last inequality being a consequence of the regular variation of  $b(t)$ . Since  $\delta$  may be chosen arbitrarily small, (2.1) follows with  $B=1$ . Thus, (3.2) holds and (3.3) is immediate.

To prove (3.4) it remains to show that

$$(3.6) \quad \lim_{t \rightarrow \infty} \frac{\log_2 b(t)}{\log_2 t} = 1.$$

However, since  $b(t)$  varies regularly at infinity, we have  $b(t) = t^e \cdot L(t)$ , where  $L(t)$  is slowly varying at infinity and thus, for an arbitrary  $\delta > 0$  and large  $t$ ,  $t^{-\delta} < L(t) < t^\delta$ , cf. [14], page 18.

Let  $\delta < \varrho$ . It follows that

$$t^{e-\delta} < b(t) < t^{e+\delta} \quad \text{for large } t,$$

from which (3.6) is immediate. Thus (3.4) holds and so does (3.5), which completes the proof.

REMARK 3.3. If  $\varrho=0$ , then (3.2) and (3.3) remain true provided  $b(t) \nearrow \infty$  as  $t \rightarrow \infty$ . As for (3.6), however, we can only assert that the limit is *at most* equal to one. If, for example  $b(t) = \log t$ , the limit equals 0, so (3.4) (and (3.5)) do not hold.

#### 4. Two applications

The first application deals with extended (non-linear) renewal theory.

Let  $\{X_n\}_{n=1}^\infty$  be i.i.d. random variables with positive expectation  $\mu$ , and positive, finite variance  $\sigma^2$ , and let  $\{S_n\}_{n=1}^\infty$  denote the partial sums. Put

$$(4.1) \quad N(t) = \min \{n; S_n > ta(n)\}, \quad t \geq 0,$$

where  $a(y)$  is a positive, ultimately non-decreasing, concave and differentiable function that varies regularly at infinity with exponent  $\alpha$ ,  $0 \leq \alpha < 1$  (cf. [7], Section 3).

Let  $\lambda = \lambda(t)$  denote the solution to the equation  $ta(y) = \mu y$  (cf. [15]), a solution which is unique for  $t$  large.

It was shown in [7], Theorem 3.5, that  $N(t)$  is asymptotically normally distributed as  $t \rightarrow \infty$ , with asymptotic mean  $\lambda(t)$  and asymptotic variance  $\frac{\sigma^2 \lambda(t)}{\mu^2 (1-\alpha)^2}$ .

The method of proof was to use Anscombe's theorem and Taylor expansion. Our first theorem will be a law of the iterated logarithm for  $\{N(t); t > 0\}$  with Theorem 4 playing the role of Anscombe's theorem.

THEOREM 5. *Under the above assumptions*

$$(4.2) \quad C \left( \left| \frac{N(t) - \lambda(t)}{\sqrt{2 \frac{\sigma^2}{\mu^2 (1-\alpha)^2} \lambda(t) \log_2 t}}; t > 0 \right| \right) = [-1, 1] \quad \text{a.s.}$$

In particular,

$$(4.3) \quad \limsup_{t \rightarrow \infty} \frac{N(t) - \lambda(t)}{\sqrt{2\lambda(t) \log_2 t}} = \frac{\sigma}{\mu(1-\alpha)} \quad \text{a.s.}$$

COROLLARY. If  $a(y) \equiv 1$ , i.e.  $N(t) = \min \{n; S_n > t\}$ , then

$$(4.4) \quad C \left( \left\{ \frac{N(t) - t/\mu}{\sqrt{2 \frac{\sigma^2}{\mu^3} t \log_2 t}}; t > 0 \right\} \right) = [-1, 1] \quad \text{a.s.}$$

In particular,

$$(4.5) \quad \limsup_{t \rightarrow \infty} \frac{N(t) - t/\mu}{\sqrt{2t \log_2 t}} = \sigma/\mu^{3/2} \quad \text{a.s.}$$

REMARK 4.1. For the case  $a(y) \equiv 1$ , Strassen versions of Theorem 5 have been proved in [16] and, provided  $E|X_1|^p < \infty$  for some  $p > 2$ , in Horváth [12]. In Horváth [10], various functional laws of the iterated logarithm for the case  $a(y) = y^\alpha$ ,  $0 \leq \alpha < 1$  are proved by strong approximation methods. Further, (4.3) was established for this case in [3], Theorem 3.3, by considering a law of the iterated logarithm for  $\max_{j \leq n} j^{-\alpha} S_j$  (see their formula (2.24)), together with inversion. (Conversely, by inverting our relation (4.3) one can reprove their formula (2.24) (in a slightly more general setting) and by inverting (4.2) the generalization to cluster sets can be obtained.)

PROOF. By definition,  $t \cdot a(\lambda(t)) = \mu\lambda(t)$ . Let  $\lambda^*(t)$  be the inverse of  $\lambda(t)$ , i.e.  $\lambda^*(t) = \mu t/a(t)$ . It follows from [9], page 22, [14], Section 1.5 that

$$(4.6) \quad \lambda(t) \text{ is regularly varying at infinity with exponent } (1-\alpha)^{-1} \geq 1.$$

Since

$$(4.7) \quad (\lambda(t))^{-1} N(t) \xrightarrow{\text{a.s.}} 1 \quad \text{as } t \rightarrow \infty$$

(see [7], Theorem 3.3) it follows that (3.1) holds (with  $\tau(t) \leftrightarrow N(t)$ ,  $b(t) \leftrightarrow \lambda(t)$ ,  $\theta=1$  and  $\varrho=(1-\alpha)^{-1} \geq 1$ ). From (3.4) we thus conclude that

$$(4.8) \quad C \left( \left\{ \frac{S_{N(t)} - \mu N(t)}{\sqrt{2\sigma^2 \lambda(t) \log_2 t}}; t > 0 \right\} \right) = [-1, 1] \quad \text{a.s.}$$

Now, since  $EX_1^2 < \infty$  is equivalent to  $n^{-1/2} X_n \xrightarrow{\text{a.s.}} 0$  as  $n \rightarrow \infty$  (cf. Example 2 in Chow and Teicher [4], page 90), it follows, in view of (4.7) that  $(\lambda(t))^{-1/2} X_{N(t)} \xrightarrow{\text{a.s.}} 0$  as  $t \rightarrow \infty$ , and in particular that

$$(4.9) \quad (2\sigma^2 \lambda(t) \log_2 t)^{-1/2} X_{N(t)} \xrightarrow{\text{a.s.}} 0 \quad \text{as } t \rightarrow \infty.$$

This fact, together with (4.8) and the relation

$$ta(N(t)) < S_{N(t)} \leq ta(N(t)) + X_{N(t)},$$

yields

$$C \left( \left\{ \frac{ta(N(t)) - \mu N(t)}{\sqrt{2\sigma^2 \lambda(t) \log_2 t}}; t > 0 \right\} \right) = [-1, 1] \quad \text{a.s.}$$

or, equivalently,

$$(4.10) \quad C\left(\left\{\frac{N(t)-(t/\mu)a(N(t))}{\sqrt{2(\sigma/\mu)^2\lambda(t)\log_2 t}}; t > 0\right\}\right) = [-1, 1] \quad \text{a.s.}$$

Suppose first that  $a(y) \equiv 1$ , i.e. that  $N(t) = \min\{n; S_n > t\}$ . Then (4.10) is the same as (4.4) and so the proof of that case is complete.

The remainder of the proof is very much like the corresponding part in [7], Theorem 3.5, and will therefore only be hinted at.

By Taylor expansion of  $a(N(t))$  around  $a(\lambda(t))$  we find, after some elementary computations, that

$$(4.11) \quad N(t) - (t/\mu)a(N(t)) = (N(t) - \lambda(t))(1 - \alpha)Y_{N(t)},$$

where

$$Y_{N(t)} = (1 - \lambda(t)a'(\lambda(t) + \varrho(N(t) - \lambda(t)))/a(\lambda(t)))/(1 - \alpha)$$

and where  $0 \leq \varrho \leq \varrho(N(t)) \leq 1$  (cf. [7], page 298). Furthermore, by [7] Lemma 3.3,

$$(4.12) \quad Y_{N(t)} \xrightarrow{\text{a.s.}} 1 \quad \text{as } t \rightarrow \infty.$$

By combining (4.10), (4.11) and (4.12) we now obtain (4.2) and the proof is complete.

The second application treats another generalization of renewal theory.

Consider a two-dimensional summation process  $\{(U_n, V_n)\}_{n=1}^\infty$  with i.i.d. summands  $\{(W_n, Z_n)\}_{n=1}^\infty$ , such that  $0 < m_w = \mathbb{E}W_1 < \infty$ ,  $0 < \sigma_w^2 = \text{Var } W_1 < \infty$  and  $0 < \sigma_z^2 = \text{Var } Z_1 < \infty$ . Further, let

$$\tau(t) = \min\{n; U_n > t\}, \quad t \geq 0.$$

The object of interest is  $V_{\tau(t)}$ .

This model was studied in [8] and it was, for example, proved that  $V_{\tau(t)}$  is asymptotically normally distributed as  $t \rightarrow \infty$  with asymptotic mean  $m_w^{-1}m_z t$  and asymptotic variance  $m_w^{-3}\gamma^2 t$ , where  $\gamma^2 = \text{Var}(m_z W_1 - m_w Z_1)$  was assumed to be positive.

By, in essence, the same proof as the proof of Theorem 5, one can prove the following law of the iterated logarithm for the stopped sums  $\{V_{\tau(t)}; t > 0\}$ .

**THEOREM 6.** *In the above setting we have*

$$(4.13) \quad C\left(\left\{\frac{V_{\tau(t)} - (m_z/m_w)t}{\sqrt{2m_w^{-3}\gamma^2 t \log_2 t}}; t > 0\right\}\right) = [-1, 1] \quad \text{a.s.}$$

*In particular*

$$(4.14) \quad \limsup_{t \rightarrow \infty} \frac{V_{\tau(t)} - m_z t/m_w}{\sqrt{2t \log_2 t}} = \gamma/m_w^{3/2} \quad \text{a.s.}$$

Note also that by choosing  $Z_n = 1$  a.s. for all  $n$  we have  $V_{\tau(t)} = \tau(t)$ ,  $\gamma^2 = \sigma_w^2$  and the above Corollary is rediscovered.

**REMARK 4.1.** Under the assumption that  $\mathbb{E}|W_1|^p < \infty$  and  $\mathbb{E}|Z_1|^p < \infty$  for some  $p > 2$  Horváth [11] proves a functional law of the iterated logarithm and a Chung-type of law of the iterated logarithm by using strong approximation methods.

(SKETCH OF) PROOF. By applying Theorem 4 to the summands  $m_w Z_n - m_z W_n$ ,  $n \geq 1$ , (cf. [8]), which have mean 0 and variance  $\gamma^2$  and the fact that  $t^{-1} \tau(t) \xrightarrow{\text{a.s.}} m_w^{-1}$  as  $t \rightarrow \infty$  (which is a consequence of extended renewal theory (cf. also (4.7) above with  $a(\cdot) \equiv 1$ )), we obtain

$$(4.15) \quad C\left(\left\{\frac{m_w V_{\tau(t)} - m_z U_{\tau(t)}}{\sqrt{2m_w^{-1}\gamma^2 t \log_2 t}}; t > 0\right\}\right) = [-1, 1] \quad \text{a.s.}$$

Further, the arguments which led to (4.9) applied to  $\{W_n\}_{n=1}^\infty$  show that  $(t \log_2 t)^{-1/2} W_{\tau(t)} \xrightarrow{\text{a.s.}} 0$  as  $t \rightarrow \infty$ , which together with (4.15) and the fact that  $0 < U_{\tau(t)} - t \leq W_{\tau(t)}$  yields the desired conclusion.

As a final remark we mention that Section 5 of [8] contains several examples in which the random quantity of interest can be described by a  $V_{\tau(t)}$ .

#### REFERENCES

- [1] ANSCOMBE, F. J., Large-sample theory of sequential estimation, *Proc. Cambridge Philos. Soc.* **48** (1952), 600—607. *MR* **14**—487.
- [2] CHANG, I. S. and HSIUNG, C. A., Strassen's invariance principle for random subsequences. *Z. Wahrsch. Verw. Gebiete* **64** (1983), 401—409.
- [3] CHOW, Y. S. and HSIUNG, A., Limiting behaviour of  $\max_{j \leq n} S_j j^{-\alpha}$  and the first passage times in a random walk with positive drift, *Bull. Inst. Math. Acad. Sinica* **4** (1976), 35—44. *MR* **53** # 11715.
- [4] CHOW, Y. S. and TEICHER, H., *Probability theory*. Springer-Verlag, New York—Heidelberg, 1978. *MR* **80a**: 60004.
- [5] CHOW, Y. S., TEICHER, H., WEI, C. Z. and YU, K. F., Iterated logarithm laws with random subsequences. *Z. Wahrsch. Verw. Gebiete* **57** (1981), 235—251. *MR* **82k**: 60062.
- [6] FELLER, W., *An introduction to probability theory and its applications*, vol. 2. John Wiley & Sons Inc. New York—London—Sydney, 1966. *MR* **35** # 1048.
- [7] GUT, A., On the moments and limit distributions of some first passage times. *Ann. Probability* **2** (1974), 277—308. *MR* **52** # 15656.
- [8] GUT, A. and JANSON, S., The limiting behaviour of certain stopped sums and some applications, *Scand. J. Statist.* **10** (1983), 281—292.
- [9] DE HAAN, L., *On regular variation and its application to the weak convergence of sample extremes*. Mathematical Centre Tracts 32, Mathematisch Centrum Amsterdam 1970. *MR* **44** # 3370.
- [10] HORVÁTH, L., Strong approximation of extended renewal processes, *Ann. Probability* **12** (1984), 1149—1166.
- [11] HORVÁTH, L., Strong approximation of certain stopped sums, *Statist. Probab. Lett.* **2** (1984), 181—185.
- [12] HORVÁTH, L., Strong approximation of renewal processes, *Stochastic Process Appl.* **18** (1984), 127—138.
- [13] HUGGINS, R. M., Laws of the iterated logarithm for time changed Brownian motion and martingales with an application to branching processes, *Ann. Probability* **13** (1985), 1148—1156.
- [14] SENETA, E., *Regularly varying functions*. Lecture Notes in Mathematics Vol. 508, Springer-Verlag, Berlin—New York, 1976. *MR* **56** # 12189.
- [15] SIEGMUND, D. O., Some one-sided stopping rules. *Ann. Math. Statist.* **38** (1967), 1641—1646. *MR* **36** # 3462.
- [16] VERVAAT, W., Functional central limit theorems for processes with positive drift and their inverses. *Z. Wahrsch. Verw. Gebiete* **23** (1972), 245—253. *MR* **47** # 9697.

(Received October 17, 1983; revised December 9, 1983)

DEPARTMENT OF MATHEMATICS  
UPPSALA UNIVERSITY  
THUNBERGSVÄGEN 3  
S-752 38 UPPSALA  
SWEDEN





# LATTICE-TILING BY CERTAIN STAR BODIES

S. STEIN

We shall investigate a problem in finite abelian groups that grows out of a question about tiling  $n$ -dimensional Euclidean space  $\mathbf{R}^n$  by translates of certain star bodies, called “semicrosses”. For positive integers,  $n$  and  $k$ , a  $(k, n)$ -semicross in  $\mathbf{R}^n$  consists of  $nk + 1$  cubes from the standard division of  $\mathbf{R}^n$  into unit cubes parallel to the coordinate axes. These  $nk + 1$  cubes consist of a corner cube together with  $n$  arms of length  $k$  in the directions of the positive axes. This is the geometric question: For which values of  $k$  can a lattice of translates of the  $(k, n)$ -semicross tile  $\mathbf{R}^n$ ? The translating vectors are assumed to have only integer coordinates and to form a group (a “ $\mathbf{Z}$ -lattice” for short).

For  $n=1$ , that is, the line, there is no restriction on  $k$  since the  $(k, 1)$ -semicross is simply an interval of length  $k + 1$ . Similarly, for  $n=2$ , the plane, there is no bound on  $k$ , since the  $(k, 2)$ -semicross, which is  $L$ -shaped, lattice-tiles  $\mathbf{R}^2$ . However, for  $n \geq 3$ ,  $k$  cannot be arbitrarily large. In [5] it was shown that if  $n \geq 3$ ,  $nk + 1$  is prime, and the  $(k, n)$ -semicross  $\mathbf{Z}$ -lattice-tiles  $\mathbf{R}^n$ , then  $k \leq 2n - 3$ . In [2] it was shown that if  $n \geq 3$  and the  $(k, n)$ -semicross  $\mathbf{Z}$ -lattice-tiles  $\mathbf{R}^n$ , then  $k \leq 2n - 1$ . However, no examples were known in which  $k$  was near these bounds. In [5] it was shown that for any odd prime  $p$  the  $(p-1, p+1)$ -semicross  $\mathbf{Z}$ -lattice tiles  $\mathbf{R}^{p+1}$ . In this case the arm length is two less than the dimension, i.e.,  $k = n - 2$ . In [1] it was shown that for  $n \geq 2$  the  $(n-1, n+1)$ -semicross  $\mathbf{Z}$ -lattice tiles  $\mathbf{R}^{n+1}$  only when  $n$  is prime. The goal of the present paper is to sharpen the bound on  $k$  to the best possible general bound, namely  $k \leq n - 2$  if the  $(k, n)$ -semicross  $\mathbf{Z}$ -lattice tiles  $\mathbf{R}^n$ ,  $n \geq 3$ . We shall deal exclusively with  $\mathbf{Z}^n$ , which serves as the discrete algebraic analog of  $\mathbf{R}^n$ .

## 1. Preliminaries

Let  $\mathbf{Z}$  denote the additive group or the ring of integers. For a positive integer  $n$  let  $\mathbf{Z}^n$  denote the additive group of  $n$ -tuples  $(x_1, x_2, \dots, x_n)$ ,  $x_i \in \mathbf{Z}$ . For positive integers  $n$  and  $k$ , the  $nk + 1$  elements of  $\mathbf{Z}^n$ ,

$$(0, 0, \dots, 0), (i, 0, \dots, 0), (0, i, 0, \dots, 0), \dots, (0, 0, \dots, 0, i),$$

$1 \leq i \leq k$ , is the  $(k, n)$ -semicross at the origin, denoted  $\mathbf{T}$ . Any translate of this set by an element  $\mathbf{v} \in \mathbf{Z}^n$ ,  $\mathbf{v} + \mathbf{T} = \{\mathbf{v} + \mathbf{t} \mid \mathbf{t} \in \mathbf{T}\}$ , is also called a  $(k, n)$ -semicross. (This

---

1980 *Mathematics Subject Classification*. Primary 52A45.

*Key words and phrases*. Tiling, lattice, abelian group.

serves as the discrete analog of the geometric semicross.) Let  $H$  be a subgroup of  $Z^n$ . If the family  $\{v+T|v\in H\}$  is a partition of  $Z^n$ , then we say that the  $(k, n)$ -semicross lattice-tiles  $Z^n$ , with lattice  $H$ . The following theorem, found in [5] or [2], reduces the question of whether such an  $H$  exists to a question on finite abelian groups.

**THEOREM 1.1.** *The  $(k, n)$ -semicross lattice-tiles  $Z^n$  if and only if there is an abelian group  $G$  of order  $nk+1$  such that in  $G$  there are  $n$  elements  $s_1, \dots, s_n$  with the property that the  $kn$  elements  $is_j$ ,  $1 \leq i \leq k$ ,  $1 \leq j \leq n$ , coincide with the non-zero elements in  $G$ .*

The subgroup  $H$  is defined as  $\{(x_1, \dots, x_n) | x_i \in Z, \sum_{i=1}^n x_i s_i = 0\}$ .

The set  $S = \{s_1, s_2, \dots, s_n\}$  is called a *splitting set* for  $S(k) = \{1, 2, \dots, k\}$ . Also,  $S(k)$  is said to *split*  $G$ .

The next theorem, proved in [4], will permit us to restrict our attention to cyclic groups.

**THEOREM 1.2.** *If  $S(k)$  splits a finite abelian group of order  $m$ , then it splits the cyclic group of order  $m$ .*

The additive cyclic group of order  $m$ ,  $Z/mZ$ , will be denoted  $C(m)$ . The same symbol will be used for  $Z/mZ$  considered as a ring.

## 2. Proof that $k \leq n-2$

We shall prove the following theorem.

**THEOREM 2.1.** *Let  $n$  and  $k$  be integers,  $n \geq 3$  and  $k \geq n-1$ . Then  $S(k)$  does not split the cyclic group  $C(nk+1)$ .*

The proof rests on two lemmas.

**LEMMA 2.2.** *Let  $n$  and  $k$  be integers,  $n \geq 3$  and  $k \geq n-1$ . Assume that  $S(k)$  splits  $C(nk+1)$ . Let  $s$  and  $s'$  be distinct elements in the splitting set, with  $s$  a generator of  $C(nk+1)$ . Then at least one of these two conditions holds:*

- (a) *There are integers  $x$  and  $y$ ,  $1 \leq x \leq n-2$ ,  $1 \leq y \leq k$ , such that  $xs + ys' = 0$ .*
- (b)  *$s' = (1-n)s$ .*

**PROOF.** Consider the  $(n-1)(k+1)$  elements

$$(1) \quad is + js',$$

$0 \leq i \leq n-2$ ,  $0 \leq j \leq k$ . If two of these are equal, there are distinct ordered pairs,  $(i, j)$  and  $(\bar{i}, \bar{j})$ , such that

$$is + js' = \bar{i}s + \bar{j}s',$$

and  $0 \leq \bar{i} \leq n-2$ ,  $0 \leq \bar{j} \leq k$ . Without loss of generality, assume that  $i \geq \bar{i}$ . Then we have

$$(2) \quad (i - \bar{i})s + (j - \bar{j})s' = 0.$$

If  $i = \bar{i}$ , then  $j \neq \bar{j}$ , and we have  $(j - \bar{j})s' = 0$ . This is a contradiction, since  $s'$  belongs to a splitting set of  $S(k)$ .

So we may assume that  $i > i$  and  $j \neq j$ . If  $J > j$ , we have  $(i-i)s = (J-j)s'$ , contradicting the fact that  $s$  and  $s'$  belong to a splitting set for  $S(k)$ . Thus  $J < j$ ; and thus (2) implies that (a) holds.

Next assume that the elements (1) are distinct. For convenience, write  $k = n-2+u$ , where  $u \geq 1$ . There are then  $(n-1)(k+1) = (n-1)(n-1+u)$  elements of the form (1) in a group of order  $nk+1 = n(n-2+u)+1$ . Hence there are exactly  $n(n-2+u)+1 - (n-1)(n-1+u) = u$  elements in  $C(nk+1)$  not of the form (1). We will show that these  $u$  elements are

$$(n-1)s, ns, \dots, (n-2+u)s = ks.$$

To show this, assume that

$$ts = is + ji',$$

where  $n-1 \leq t \leq k$ ,  $0 \leq i \leq n-2$ ,  $0 \leq j \leq k$ . Subtraction yields

$$(t-i)s = js',$$

with  $1 \leq t-i \leq k$ ,  $0 \leq j \leq k$ , contradicting the assumption that  $s$  and  $s'$  belong to a splitting set for  $S(k)$ .

Next consider the representation of the remaining elements,  $0, s, \dots, (n-2)s, (k+1)s, (k+2)s, \dots, nks$  in the form (1). That is, consider the representation of  $\{ws \mid 0 \leq w \leq n-2 \text{ or } k+1 \leq w \leq kn\}$ .

For  $0 \leq w \leq n-2$ , we have the representation  $ws = ws + 0s'$ .

Consider next the representation of  $(k+1)s$  in the form (1),

$$(k+1)s = is + js',$$

$0 \leq i \leq n-2$ ,  $0 \leq j \leq k$ . Then

$$(3) \quad (k+1-i)s = js'.$$

If  $i > 0$ , (3) violates the fact that  $s$  and  $s'$  belong to a splitting set for  $S(k)$ . Thus  $i = 0$ , and we have

$$(k+1)s = js',$$

for some  $j$ ,  $1 \leq j \leq k$ .

Before examining the representation of further elements,  $ws$ , we denote the  $j$  just obtained by  $j_0$ . Thus  $(k+1)s = j_0s'$ .

It follows that the representation of the elements of the form  $(k+1)s + xs$ ,  $1 \leq x \leq n-2$ , in the form (1) is  $xs + j_0s'$ .

Consider the representation of  $(k+1)s + (n-2)s + s = (k+n)s$  in the form (1):

$$(k+n)s = is + js',$$

$0 \leq i \leq n-2$ ,  $0 \leq j \leq k$ . If  $i > 0$ , we have

$$(4) \quad [k+(n-i)]s = js'.$$

Combined with the inequality

$$k+2 \leq k+n-i \leq k+n-1,$$

(4) contradicts the representation of  $[k+(n-i)]s$  already obtained. Thus  $i = 0$  and

we may write

$$(k+n)s = j_1 s,$$

$$1 \leq j_1 \leq k.$$

In a similar manner, we find that

$$[k+1+2(n-1)]s = j_2 s',$$

then that

$$[k+1+3(n-1)]s = j_3 s',$$

and, more generally,

$$[k+1+v(n-1)]s = j_v s',$$

for  $0 \leq v \leq k-1$ . In each case,  $1 \leq j_v \leq k$ .

Let us compare  $j_v$  and  $j_{v+1}$  for each  $v$ ,  $0 \leq v \leq k-2$ . We have

$$[k+1+v(n-1)]s = j_v s'$$

and

$$[k+1+(v+1)(n-1)]s = j_{v+1} s'.$$

Subtraction yields

$$(n-1)s = (j_{v+1} - j_v)s'.$$

Since  $s$  and  $s'$  belong to a splitting set for  $S(k)$ ,  $j_{v+1} - j_v < 0$ , hence  $j_{v+1} < j_v$ . Thus the  $k$  integers  $j_0, j_1, \dots, j_{k-1}$  are a descending sequence in the interval  $[1, k]$ . Therefore, they must be, in order, the integers  $k, k-1, \dots, 1$  respectively. In particular  $j_{k-1} = 1$ . This means that

$$[k+1+(k-1)(n-1)]s = 1s'.$$

Since this equation holds in the cyclic group  $C(nk+1)$ , where  $nk+1=0$ , it is equivalent to the equation  $(1-n)s = s'$ . This establishes (b) and completes the proof of the lemma.

Incidentally, (a) and (b) in Lemma 2.2 are exclusive.

**LEMMA 2.3.** *Let  $n$  and  $k$  be integers,  $n \geq 3$  and  $k \geq n-1$ . Assume that  $S(k)$  splits  $C(nk+1)$ . Let  $s$  and  $s'$  be distinct elements in the splitting set, with  $s$  not a generator of  $C(nk+1)$ . Then there are integers  $x$  and  $y$ ,  $1 \leq x \leq n-2$ ,  $1 \leq y \leq k$ , such that  $xs + ys' = 0$ .*

**PROOF.** As in the proof of Lemma 2.2, consider the  $(n-1)(k+1)$  elements (1). If two of these are equal, the lemma follows, as in the proof of Lemma 2.2. We will show that the case where all the elements (1) are distinct cannot occur.

Write  $s = s^*d$ , where  $s^* = \gcd(s, nk+1)$ , hence  $\gcd(d, nk+1) = 1$ . There is an element  $d_1$  in  $C(nk+1)$  such that  $dd_1 = 1$ . Replace the entire splitting set,  $s = s_1, s' = s_2, \dots, s_n$  by  $sd_1, s'd_1, \dots, s_nd_1$ . Observe that  $sd_1 = s^*$ , a divisor of  $nk+1$ . Moreover, the proof of the lemma for the new splitting set will imply the truth of the lemma for the original splitting set.

So, without loss of generality, we may begin, with the pair  $s$  and  $s'$  from a splitting set such that  $s$  divides  $nk+1$ .

Consider the additive subgroup of  $C(nk+1)$  generated by  $s$ , which consists of the  $(nk+1)/s$  elements  $0, s, 2s, \dots, nk+1-s$ . All the elements

$$(5) \quad is + js',$$

$0 \leq i \leq n-2$ ,  $0 \leq j \leq k$ ,  $j$  a multiple of  $s$ , lie in this subgroup. Moreover, as in the argument for Lemma 2.2, none of the elements (5) can coincide with any of the elements  $(n-1)s, ns, \dots, ks$ . (There are  $k-n+2$  such excluded elements.)

We will show that there are more elements of the form (5) than there are in the set obtained by deleting  $(n-1)s, ns, \dots, ks$  from the subgroup generated by  $s$ .

Let  $k=qs+r$ ,  $0 \leq r \leq s-1$ . The number of elements of the form (5) is  $(n-1)(q+1)$ , which equals

$$(n-1) \left( \frac{k-r}{s} + 1 \right).$$

The desired inequality therefore is

$$(n-1) \left( \frac{k-r}{s} + 1 \right) > \frac{kn+1}{s} - (k-n+2),$$

which is equivalent, by straightforward algebraic manipulation, to the inequality

$$r < \frac{k+1}{n-1} (s-1).$$

But this last inequality holds since  $r \leq s-1$  and  $k+1 > n-1$ . This concludes the proof of the lemma.

With the two lemmas available, the proof of Theorem 2.1 is short.

**PROOF OF THEOREM 2.1.** Assume, without loss of generality that 1 is in the splitting set. Let  $s_2, s_3, \dots, s_n$  be the other elements in the splitting set. For each index  $p$ ,  $2 \leq p \leq n$ , consider the pair of elements 1 and  $s_p$ . According to Lemma 2.2, either there are integers  $x_p$  and  $y_p$ ,  $1 \leq x_p \leq n-2$ ,  $1 \leq y_p \leq k$  such that  $x_p + y_p s_p = 0$  or else  $s_p = (1-n)1 = 1-n$ . If the first case (a) holds for each  $p$ ,  $2 \leq p \leq n$ , then there would be  $p-1$  values  $x_p$  in the interval  $[1, n-2]$ . By the pigeon-hole principle, two of them are equal, say  $x_{p_1} = x_{p_2}$ . From the equations  $x_{p_1} + y_{p_1} s_{p_1} = 0$ ,  $x_{p_2} + y_{p_2} s_{p_2} = 0$ , and  $x_{p_1} = x_{p_2}$ , it follows that  $y_{p_1} s_{p_1} = y_{p_2} s_{p_2}$ , violating the assumption that  $s_{p_1}$  and  $s_{p_2}$  belong to a splitting set for  $S(k)$ . Thus there is an index  $p$  such that  $s_p = (1-n)1 = 1-n$ .

So now we have  $1-n$  in the splitting set. There are two cases to consider:  $1-n$  is a generator of  $C(nk+1)$  or  $1-n$  is not a generator of  $C(nk+1)$ .

In the first case, reasoning similar to that just encountered shows that  $(1-n)^2$  is an element of the splitting set. This element is different from 1 and  $1-n$ . Let  $m=nk+1$ . We have  $m=nk+1=n(n-2+u)+1=(1-n)^2+un$ . Thus  $(1-n)^2 \equiv -un \pmod{m}$ . Hence

$$k((1-n)^2) \equiv -kun \equiv u \pmod{m}.$$

Since  $0 \leq u \leq k$  (in fact  $1 \leq u \leq k-1$ ), this equation implies that

$$k((1-n)^2) \equiv u \cdot 1 \pmod{m},$$

violating the fact that 1 and  $(1-n)^2$  belong to the splitting set.

Finally, we consider the case where  $1-n$  is not a generator of  $C(nk+1)$ , that is,  $\gcd(1-n, nk+1) > 1$ . In this case, apply Lemma 2.3 with  $s=1-n$  and  $s'=s_2, s_3, \dots, s_n$  successively. For each index  $p$ ,  $2 \leq p \leq n$  there are integers  $x_p$ ,

and  $y_p$  such that  $x_p(1-n) + y_p s_p = 0$ ,  $1 \leq x_p \leq n-2$ ,  $1 \leq y_p \leq k$ . As earlier in this proof, we again obtain a contradiction. This proves the theorem.

**COROLLARY 2.4.** *Let  $n$  be one more than a composite number. If  $S(k)$  splits  $C(nk+1)$ , then  $k \leq n-3$ .*

This follows from Theorem 2.1 and the theorem in [1] quoted in the introduction.

Theorem 2.1 sheds light on the question, "For which values of  $n$  and  $k$  does  $S(k)$  split  $C(nk+1)$ ?" For other work on this and related questions see [1]–[9].

**ADDED IN PROOF.** Using different methods, S. Szabó has obtained a slightly weaker bound for  $k$  in semicrosses but also the strongest bound yet found in the case of the  $(k, n)$ -cross. See S. Szabó, A bound for  $k$  for tiling by  $(k, n)$ -crosses and semicrosses, *Acta Math. Hungar.* **44** (1984), 97–99.

#### REFERENCES

- [1] GALOVICH, S. and STEIN, S., Splitting of abelian groups by integers, *Aequationes Math.* **22** (1981), 249–267. *MR* **83g**: 20018.
- [2] HAMAKER, W., Factoring groups and tiling space, *Aequationes Math.* **9** (1973), 145–149. *MR* **48** # 5893.
- [3] HAMAKER, W. and STEIN, S., Combinatorial packing of  $\mathbb{R}^3$  by certain error spheres, *IEEE Trans. Information Theory*.
- [4] HICKERSON, D., Splittings of finite groups, *Pacific J. Math.*, **107** (1983), 141–171.
- [5] STEIN, S., Factoring by subsets, *Pacific J. Math.* **22** (1967), 523–541. *MR* **36** # 2517.
- [6] STEIN, S., Packings of  $\mathbb{R}^n$  by certain error spheres, *IEEE Trans. Information Theory* **IT**—**30** (1984), no. 2, 356–363.
- [7] STEIN, S., Tiling, packing, and covering by clusters, *Rocky Mountain J. of Math.* **16** (1986), 277–321.
- [8] SZABÓ, S., Rational tilings by  $n$ -dimensional crosses I, *Proc. Amer. Math. Soc.* **87** (1983), 213–222.
- [9] SZABÓ, S., Rational tilings by  $n$ -dimensional crosses II.

(Received November 9, 1983)

DEPARTMENT OF MATHEMATICS  
UNIVERSITY OF CALIFORNIA AT DAVIS  
DAVIS, CA 95616  
U.S.A.



# THE DISTRIBUTION OF VALUES OF A CLASS OF ARITHMETIC FUNCTIONS

M. V. SUBBARAO\* and R. SITARAMACHANDRARAO\*\*

## Abstract

The distribution of values taken by the Euler totient function and the sum of the divisors function have been investigated by several authors. In this paper, we prove an asymptotic formula, with an estimate for the error term, for the number of integers  $n$  satisfying  $1 \leq f(n) \leq x$  where  $f$  belong to a prescribed class of nonnegative integer valued multiplicative function satisfying some very general conditions. From our main theorem, we deduce as special cases earlier results concerning the above mentioned functions due to P. Erdős, P. T. Bateman and Ivić. Besides, we give several other applications which are believed to be new, some involving arithmetical functions not considered before.

## 1. Introduction

The distribution of values taken by  $\varphi$ , the Euler totient function, has been investigated from many points of view. For example, more than 50 years ago, T. Vijayaraghavan [24] proved that the sequence  $\left\{ \frac{\varphi(n)}{n} \right\}_{n=1,2,\dots}$  is everywhere dense in the unit interval while a classic result due to H. Weyl implies that this is not uniformly distributed (mod 1). Another classical result is due to I. J. Schoenberg [20] who proved that this sequence has a continuous distribution function. That is, there exists a continuous monotonic function  $a$  with  $a(0)=0$ ,  $a(1)=1$  satisfying for each  $\alpha \in [0, 1]$

$$\frac{1}{x} \# \{n \in [1, x] | \varphi(n)/n \leq \alpha\} \rightarrow a(\alpha)$$

as  $x \rightarrow \infty$ . In 1945, P. Erdős [12] took a different approach in that he studied  $A(x)$ , the number of positive integers  $n$  satisfying  $\varphi(n) \leq x$ . In particular, he proved that  $\lim_{x \rightarrow \infty} A(x)/x$  exists. While Erdős' proof rests on Schoenberg's work, a completely elementary proof of it was given by R. E. Dressler [9] who also evaluated the limit. Once the existence of the limit is known, it is comparatively easy to evaluate it, for example by an abelian argument, as noted by P. T. Bateman [1]. Thus we have the relation

$$(1.1) \quad \lim_{x \rightarrow \infty} \frac{A(x)}{x} \frac{\zeta(2)\zeta(3)}{\zeta(6)} = \alpha,$$

\* This research is supported in part by a Canadian National Research Council Grant.

\*\* Presently at the University of Alberta, Edmonton, on leave from Andhra University.

1980 *Mathematics Subject Classification*. Primary 10H25.

*Key words and phrases*. Euler totient function, sum of the divisors function, Riemann Zeta function and generalized integers.

say, where  $\zeta$  denotes the Riemann Zeta function. In 1972, P. T. Bateman [1] illustrated, in an interesting way, three techniques in analytic number theory to obtain estimates for the error term  $A(x) - \alpha x$  where  $\alpha$  is given by (1.1). For example, he proved that for fixed  $\varepsilon > 0$  and as  $x \rightarrow \infty$

$$(1.2) \quad A(x) = \alpha x + O\left\{x \exp\left\{- (1 - \varepsilon) \frac{1}{2} \log x (\log \log x)^{1/2}\right\}\right\}.$$

He conjectured that  $A(x) = \alpha x + O(x \exp\{-(\log x)^{1-\varepsilon}\})$  holds for every positive  $\varepsilon$  and that  $A(x) - \alpha x \neq o(x^\lambda)$  for any  $\lambda < 1$ . As far as we can determine these are still open. At the end of his paper, Bateman also noted that the various methods developed in his paper in the case of  $\varphi(n)$  would also work for  $\sigma(n)$ , the sum of all the divisors of  $n$ , and, in particular, a result corresponding to (1.2) above could be obtained.

Recently, A. Ivić (cf. [15], Theorem 2) obtained a result similar to (1.2) above for a class of multiplicative functions. In fact, he considered functions  $f$  which are of the form: For primes  $p$  and positive integral  $m$

$$(1.3) \quad f(p^m) = p^m + a_{1,m} p^{m-1} + \dots + a_{m,m}$$

where  $-1 \leq a_{i,m} \leq k$  for some  $k \geq 0$  and all  $i = 1, 2, \dots, m$ . His method of proof is essentially the same as that of P. T. Bateman [1]. However, the order estimate for the error term is better than that of P. T. Bateman (from method A). Since Ivić used Walfisz's estimate for the error term in prime number theorem rather than that of de la Vallée Poussin.

The class of functions considered by Ivić does not contain some of the well-known functions, namely,  $\sigma_r(n)$  (the sum of the  $r$ th powers of the divisors of  $n$ ),  $J_k(n)$  (the Jordan totient function of order  $k$ ) and  $\Phi_n(n)$  (the Schemmel's totient function of order  $k$ ). This can be said to be one of the limitations of Ivić's paper.

The main object of this paper is to obtain a fairly general theorem for a wide class of multiplicative functions which includes all the functions stated by Ivić [15] in his introduction and moreover several others. It may be noted that the sets of functions considered by Ivić and us contain several arithmetic functions in common but neither is a subset of the other. If, however, we restrict the coefficients  $a_{i,m}$  in (1.3) to satisfy  $-c \leq a_{i,m} \leq k$  for  $m \geq 2$  where  $0 < c < 1$  and  $k \geq 0$ , then our results are applicable to such functions. Section 4 contains a rich class of illustrations which result from an application of our Theorem 2.1 to the well-known functions  $J_k$  (the Jordan totient function of order  $k$ ),  $\sigma_k$  (sum of the  $k$ th powers of all the positive divisors), their unitary analogues, Schemmel's totient function of order  $k$  and Dedekind's  $\psi$ -function. In Section 5, an  $M$ -void analogue of the Euler totient function, which does not seem to have been mentioned in the literature, is introduced which affords us with another class of illustrations of Theorem 2.1. In Section 6, the scope of applicability of Theorem 2.1 is extended.

## 2. Notation and statement

Let  $\mathbf{P}$  denote the set of all prime numbers and  $\mathbf{Z}^+$ , the set of all positive integers. For a subset  $A$  of  $\mathbf{P}$ , we write  $S(A)$  to mean the multiplicative semigroup generated by  $A$ . Let  $\mathcal{F}$  denote the class of all integer valued multiplicative arithmetic functions  $f$  to each of which there correspond subsets  $P_1$  and  $P_2$  satisfying the following conditions:

(2.1)  $P_1$  and  $P_2$  are finite sets, possibly empty.

(2.2)  $f(p^m)=0$  for all  $p \in P_1$  and  $m \in \mathbf{Z}^+$ .

(2.3) There exists a  $\beta > 1$  such that

$$f(p^m) \leq \beta^{m-1} \quad \text{for all } p \in P_2 \text{ and } m \in \mathbf{Z}^+.$$

Let  $P_3 = \mathbf{P} \setminus (P_1 \cup P_2) = \{p_1, p_2, p_3, \dots\}$  with  $p_1 < p_2 < p_3 < \dots$ .

(2.4) There exists a  $u \in \mathbf{Z}^+$  and an integer  $v > -2^u + 1$  such that

$$f(p) = p^u + v \quad \text{for all } p \in P_3 \text{ and}$$

(2.5) There exist constants  $B > 0$  and  $\delta > u/2$  such that

$$f(p^m) \leq Bp^{m\delta} \quad \text{for all } p \in P_3 \text{ and } m \geq 2.$$

We note that (2.3), (2.4) and (2.5) together imply that  $f(p^m) \rightarrow \infty$  as  $p^m \rightarrow \infty$ ,  $p \in P_2 \cup P_3$ . Hence  $f(n) \rightarrow \infty$  as  $n \rightarrow \infty$  through the semigroup  $S(P_2 \cup P_3)$ . Hence for each  $f \in \mathcal{F}$  and  $m \in \mathbf{Z}^+$ , the number  $a_m$  of solutions in  $n$  of the equation  $f(n) = m$  is finite. Thus on writing  $A(x; f) = \sum_{1 \leq n \leq x} a_m$ , we see that  $A(x, f)$  is the number of positive integers  $n$  with  $1 \leq f(n) \leq x$ .

The main result of the paper is the following

**THEOREM 2.1.** *Let  $f \in \mathcal{F}$ . Then as  $x \rightarrow \infty$*

$$A(x; f) = A(f)x^{1/u} + O(x^{1/u} \exp \{-c(\log x)^{3/8-\epsilon}\})$$

*for each  $\epsilon > 0$  where  $A(f)$  is a positive constant given by*

$$(2.6) \quad A(f) = \lim_{s \rightarrow 1+} \left[ \prod_{p \in \mathbf{P}} (1 - p^{-s}) \prod_{p \in P_2 \cup P_3} \left\{ \sum_{m=0}^{\infty} (f(p^m))^{-s} \right\} \right]$$

*or equivalently by*

$$A(f) = \lim_{s \rightarrow 1+} [(s-1) \prod_{p \in P_2 \cup P_3} \left\{ \sum_{m=0}^{\infty} (f(p^m))^{-s} \right\}]$$

*and  $c$  is any positive constant.*

## 3. Proof of Theorem 2.1

For the proof of the theorem, we require a result in Beurling's theory of generalized integers. This result (Lemma 1 below), in its present form, is due to H. G. Diamond [6] and is refinement over earlier results due to B. Nyman [18] and P. Malliavin [16].

Suppose that  $1 < l_1 \leq l_2 \leq l_3 \leq \dots, l_n \rightarrow \infty$  is a sequence of positive numbers. Let the product  $\prod_{n=1}^{\infty} (1 - l_n^{-s})^{-1} = \prod_{n=1}^{\infty} (1 + l_n^{-s} + l_n^{-2s} + \dots)$  be formally expanded into a General Dirichlet series  $\sum_{n=1}^{\infty} \beta_n g_n^{-s}$  where  $g_1 = 1, g_2, g_3, \dots$  is an increasing sequence of positive numbers with range containing  $S(\{l_1, l_2, \dots\})$  and  $\beta_1 = 1, \beta_2, \beta_3, \dots$  are nonnegative integers.

LEMMA 1. *With the above notation, suppose that*

$$\sum_{l_n \leq x} 1 = \int_2^x \frac{du}{\log u} + O(x \exp \{-b(\log x)^a\})$$

where  $b > 0$  and  $0 < a < 1$ . Then

$$\sum_{g_n \leq x} \beta_n = Cx + O(x \exp \{-c_2(\log x)^{a/(a+1)}\})$$

for every  $c_2 > 0$ . Here the constant  $C$  is given by

$$\begin{aligned} C &= \lim_{s \rightarrow 1+} \left[ \prod_{p \in P} (1 - p^{-s}) \prod_{n=1}^{\infty} (1 - l_n^{-s})^{-1} \right] \\ (3.1) \quad &= \lim_{s \rightarrow 1+} [(s-1) \prod_{n=1}^{\infty} (1 - l_n^{-s})^{-1}]. \end{aligned}$$

For convenience of later reference, we also state

LEMMA 2 (cf. [13], Theorem 41). *If  $f$  is multiplicative and  $\prod_{p \in P} \left\{ \sum_{m=0}^{\infty} |f(p^m)| \right\}$  converges, then  $\sum_{n=1}^{\infty} f(n)$  converges absolutely and*

$$\sum_{n=1}^{\infty} f(n) = \prod_{p \in P} \left\{ \sum_{m=0}^{\infty} f(p^m) \right\}.$$

We now proceed to the proof of the theorem. We consider the series  $\sum_{n \in S(P_2 \cap P_3)} (f(n))^{-s}$ . Assuming its absolute convergence for some  $s > 0$  (which will be justified soon), we have

$$\begin{aligned} \sum_{n \in S(P_2 \cup P_3)} (f(n))^{-s} &= \prod_{p \in P_2} \left\{ \sum_{m=0}^{\infty} (f(p^m))^{-s} \right\} \prod_{p \in P_3} \left\{ \sum_{m=0}^{\infty} (f(p^m))^{-s} \right\} = \\ (3.2) \quad &= \prod_{p \in P_2} \{1 - (f(p))^{-s}\}^{-1} \left[ \prod_{p \in P_2} \left\{ \sum_{m=0}^{\infty} (f(p^m))^{-s} \right\} \prod_{p \in P_3} \{1 - (f(p))^{-2s} + \right. \\ &\quad \left. + (1 - (f(p))^{-s}) \sum_{m=2}^{\infty} (f(p^m))^{-s} \right] \\ &= \Pi_1(s) \Pi_2(s), \end{aligned}$$

say. Further, since  $f$  is positive integer valued on  $S(P_2 \cup P_3)$ , we write  $\Pi_1(s) = \sum_{n=1}^{\infty} \frac{b_n}{n^s}$

and  $\Pi_2(s) = \sum_{n=1}^{\infty} \frac{d_n}{n^s}$ .  $P_2$  being a finite set, the product over  $P_2$  in the above converges absolutely for each  $s > 0$ . Also concerning the product over  $P_3$  in the extreme right, we note that, by (2.4) and (2.5), for  $s > 0$  and as  $p \rightarrow \infty$

$$(f(p))^{-2s} = (p^u + v)^{-2s} \sim p^{-2us},$$

$$(f(p))^{-s} \rightarrow 0$$

and

$$\sum_{m=2}^{\infty} (f(p^m))^{-s} \ll \sum_{m=2}^{\infty} p^{-m\delta s} = p^{-2\delta s} / (1 - p^{-\delta s}) \ll p^{-2\delta s}.$$

Consequently, the product defining  $\Pi_2$  and hence by Lemma 2 the series  $\sum_{n=1}^{\infty} d_n n^{-s}$  converges absolutely for  $s > \max\left(\frac{1}{2u}, \frac{1}{2\delta}\right)$ . Further, since  $(f(p))^{-s} \sim p^{-us}$  as  $p \rightarrow \infty$ , the product defining  $\Pi_1$  converges absolutely for  $s > 1/u$ . Hence by another application of Lemma 2, we see that the series appearing on the left of (3.2) converges absolutely for  $s > 1/u$ .

Now we apply Lemma 1 with  $l_n = (p_n^u + v)^{1/u}$ ,  $g_n = n^{1/u}$  and  $\beta_n = b_n$ . For this purpose, we observe that for each  $\varepsilon > 0$

$$\sum_{\substack{n \\ (p_n^u + v)^{1/u} \leq x}} 1 = \int_2^x \frac{dt}{\log t} + O(x \exp\{-c_3(\log x^{(3/5)-\varepsilon})\})$$

obtainable from A. Walfisz [25],  $c_3$  being a positive constant. Combining this with Lemma 1, we obtain

$$B(x) = \sum_{n \leq x} b_n = Cx^{1/u} + O(x^{1/u} \exp\{-c_2(\log x)^{(3/8)-\varepsilon}\})$$

for every  $\varepsilon > 0$  and  $c_2 > 0$ . Here, by (3.1)

$$(3.3) \quad C = \lim_{s \rightarrow 1+} \left[ \prod_{p \in P} (1 - p^{-s}) \prod_{p \in P_3} \{1 - (p^u + v)^{-s-1}\} \right].$$

Since the series  $\sum_{n=1}^{\infty} d_n n^{-s}$  has abscissa of absolute convergence at most  $\max\left(\frac{1}{2u}, \frac{1}{2\delta}\right) < \frac{1}{u}$ , an elementary argument then yields

$$\begin{aligned} A(x; f) &= \sum_{1 \leq n \leq x} a_n = \sum_{n=1}^{[x]} d_n B(x/n) = \\ &= x^{1/u} C \sum_{n=1}^{\infty} d_n n^{-1/u} + O(x^{1/u} \exp\{-c(\log x)^{(3/8)-\varepsilon}\}) \end{aligned}$$

for every  $\varepsilon > 0$  and  $c > 0$ . Also by (3.3) and (2.6)

$$\begin{aligned} C \sum_{n=1}^{\infty} d_n n^{-1/u} &= \lim_{s \rightarrow 1+} \left[ \prod_{p \in P} (1-p^{-s}) \prod_{p \in P_3} \{1-(f(p))^{-s/u}\}^{-1} \right. \\ &\cdot \prod_{p \in P_2} \left\{ \sum_{m=0}^{\infty} (f(p^m))^{-s/u} \right\} \cdot \prod_{p \in P_3} \{1-(f(p))^{-2s/u} + (1-(f(p))^{-s/u}) \sum_{m=2}^{\infty} (f(p^m))^{-s/u}\} \Big] = \\ &= \lim_{s \rightarrow 1+} \left[ \prod_{p \in P} (1-p^{-s}) \prod_{p \in P_2 \cup P_3} \left\{ \sum_{m=0}^{\infty} (f(p^m))^{-s/u} \right\} \right] = A(f). \end{aligned}$$

This completes the proof of Theorem 2.1.

#### 4. Applications

In this section, we illustrate Theorem 2.1 by specializing it to various well-known arithmetic functions. Let  $K$  and  $s$  be positive integers with  $s \leq 2$ . Let  $\mu$  denote the Möbius function and  $\mu^*$  its unitary analogue defined by  $\mu^*(n) = (-1)^{\omega(n)}$  where  $\omega(n)$  denotes the number of distinct prime divisors of  $n$ . The symbol  $d \parallel n$  means that  $d \mid n$  and  $(d, n/d) = 1$ . Let the arithmetic functions  $\sigma_K$ ,  $\varphi_{K,s}$ ,  $\sigma_K^*$ ,  $\varphi_K^*$  and  $\Phi_K$  be defined by

$$\sigma_K(n) = \sum_{d \mid n} d^K, \quad \varphi_{K,s}(n) = \sum_{d \mid n} (\mu(d))^s \left(\frac{n}{d}\right)^K,$$

$$\sigma_K^*(n) = \sum_{d \parallel n} d^K, \quad \varphi_K^*(n) = \sum_{d \parallel n} \mu^*(d) \left(\frac{n}{d}\right)^K$$

and

$$\Phi_K(n) = \begin{cases} n \prod_{p \mid n} \left(1 - \frac{K}{p}\right) & \text{if each prime divisor of } n \text{ is greater than } K, \\ 0 & \text{otherwise.} \end{cases}$$

For the sake of shortness, we write  $\varphi_{K,1} = J_K$ ,  $\varphi_{K,2} = \psi_K$ ,  $\psi_1 = \psi$ ,  $\sigma_1^* = \sigma^*$  and  $\varphi_1^* = \varphi^*$ .  $J_K$  is the well-known Jordan totient function of order  $K$  (cf. [8], p. 147),  $\psi_K$  is an extension of Dedekind's  $\psi$ -function (cf. [8], p. 123) and  $\Phi_K$  is the Schemmel's totient function of order  $K$  (cf. [8], p. 147). Clearly,  $J_1 = \Phi_1 = \varphi$ , the Euler totient function.  $\sigma^*$  and  $\varphi^*$  are the respective unitary analogues of the functions  $\sigma$  and  $\varphi$  [2].

It is clear that for  $K \in \mathbb{Z}^+$ , each of the functions  $\sigma_K$ ,  $J_K$ ,  $\psi_K$ ,  $\sigma_K^*$ ,  $\varphi_K^*$  and  $\Phi_K$  belongs to  $\mathcal{F}$ . Hence Theorem 2.1 could be specialized to these functions and we obtain the following

**THEOREM 4.1.** *For each  $c > 0$ , we have, as  $x \rightarrow \infty$*

$$(4.1) \quad A(x; \sigma_K) = A(\sigma_K) x^{1/K} + O(x^{1/K} (\delta(x))^c),$$

$$(4.2) \quad A(x, J_K) = A(J_K) x^{1/K} + O(x^{1/K} (\delta(x))^c),$$

$$(4.3) \quad A(x, \psi_K) = A(\psi_K) x^{1/K} + O(x^{1/K} (\delta(x))^c),$$

$$(4.4) \quad A(x, \sigma_K^*) = A(\sigma_K^*) x^{1/K} + O(x^{1/K} (\delta(x))^c),$$

$$(4.5) \quad A(x, \varphi_K^*) = A(\varphi_K^*) x^{1/K} + O(x^{1/K} (\delta(x))^c)$$



and

$$(4.6) \quad A(x; \Phi_K) = A(\Phi_K)x + O(x(\delta(x))^c)$$

where  $c$  is any positive constant,

$$(4.7) \quad \delta(x) = \exp \{-(\log x)^{(3/8)-\epsilon}\},$$

and

$$A(\sigma_K) = \prod_{p \in \mathbf{P}} \left\{ (1-p^{-1}) \sum_{m=0}^{\infty} (1+p^K + \dots + p^{mK})^{-1/K} \right\},$$

$$A(J_K) = \prod_{p \in \mathbf{P}} \left\{ 1-p^{-1} + p^{-1}(1-p^{-K})^{-1/K} \right\},$$

$$A(\psi_K) = \prod_{p \in \mathbf{P}} \left\{ 1-p^{-1} + (p^K+1)^{-1/K} \right\},$$

$$A(\sigma_K^*) = \prod_{p \in \mathbf{P}} \left\{ (1-p^{-1}) \sum_{m=0}^{\infty} (1+p^{mK})^{-1/K} \right\},$$

$$A(\varphi_K^*) = \prod_{p \in \mathbf{P}} \left\{ (1-p^{-1}) \left( 1 + \sum_{m=1}^{\infty} (p^{mK}-1)^{-1/K} \right) \right\},$$

and

$$A(\Phi_K) = \frac{\varphi(K)}{K} \prod_{p > K} \left( 1 - \frac{1}{p} + \frac{1}{p-K} \right).$$

In particular,  $A(J_1) = A(\Phi_1) = \zeta(2)\zeta(3)/\zeta(6)$ .

REMARK 4.1. As mentioned in the introduction, (4.1) and (4.2) in case  $K=1$  in their asymptotic forms and without explicit determination of the constants  $A(\varphi)$  and  $A(\sigma)$  were due to P. Erdős [12]. R. E. Dressler [9], [10] gave completely elementary proofs of Erdős' results and also evaluated the constants. As mentioned already, P. T. Bateman [1] established (4.1) and (4.2) in case  $K=1$  with sharper estimates for the error terms. H. G. Diamond [7] discussed an extension of  $A(x; \varphi)$ . The result (4.5), in case  $K=1$ , in its asymptotic form and without an explicit determination of the constant  $A(\varphi^*)$  was due to M. Ismail and M. V. Subbarao [14].

## 5. Some more applications

Let  $M$  be a set of positive integers with  $\min M = s \geq 2$ . Following G. J. Rieger [19], we say that an integer is  $M$ -void if it is positive and in its canonical factorisation  $\prod_p p^{l_p}$ , no  $l_p$  belongs to  $M$ . Let  $\mathcal{Q}_M$  denote the set of all  $M$ -void integers and  $q_M$  its characteristic function. Let  $\lambda_M$  denote the inversion function of the set  $\mathcal{Q}_M$ , that is,  $\lambda_M$  is the unique arithmetic function defined by

$$q_M(n) = \sum_{d|n} \lambda_M(d)$$

for all  $n$ . It is clear that  $q_M$  and  $\lambda_M$  are multiplicative functions and further that for  $p \in \mathbf{P}$  and  $m \in \mathbf{Z}^+$

$$(5.1) \quad \lambda_M(p^m) = \begin{cases} -1 & \text{if } m \in M, m-1 \notin M, \\ 1 & \text{if } m \notin M, m-1 \in M, \\ 0 & \text{otherwise.} \end{cases}$$

For positive integral  $n$ , we define the  $M$ -void analogue of the Euler totient function to be the number of integers  $x$  in  $[1, n]$  such that the g.c.d.  $(x, n) \in Q_M$ . Denoting this by  $\varphi_M(n)$ , we have

$$(5.2) \quad \varphi_M(n) = \sum_{d|n} \lambda_M(d) \frac{n}{d}.$$

In fact,

$$\begin{aligned} \varphi_M(n) &= \sum_{\substack{1 \leq x \leq n \\ (x, n) \in Q_M}} 1 = \sum_{1 \leq x \leq n} q_M((x, n)) = \\ &= \sum_{1 \leq x \leq n} \sum_{d|(x, n)} \lambda_M(d) = \sum_{d|n} \lambda_M(d) \sum_{\substack{1 \leq x \leq n \\ d|x}} 1 = \sum_{d|n} \lambda_M(d) \frac{n}{d}. \end{aligned}$$

Now by (5.1) and (5.2), we have, for  $p \in \mathbf{P}$  and  $m \in \mathbf{Z}^+$

$$(5.3) \quad \varphi_M(p^m) = \begin{cases} p^m & \text{if } 1 \leq m \leq s-1, \\ p^m + \lambda_M(p^s)p^{m-s} + \dots + \lambda_M(p^m) & \text{if } m \geq s. \end{cases}$$

Thus for  $m \geq s$

$$\begin{aligned} \varphi_M(p^m) &\geq p^m \left( 1 - \frac{1}{p^s} - \frac{1}{p^{s+1}} \dots - \frac{1}{p^m} \right) \geq \\ &\geq p^m \left( 1 - \frac{1}{2^2} - \frac{1}{2^3} - \dots \right) = \frac{p^m}{2}. \end{aligned}$$

Combining this with (5.3), we conclude that for  $p \in \mathbf{P}$  and  $m \in \mathbf{Z}^+$ ,  $\varphi_M(p) = p$  and  $\varphi_M(p^m) \geq (1/2)p^m$ .

Now it is clear that  $\varphi_M \in \mathcal{F}$  and as such Theorem 2.1 yields the following

**THEOREM 5.1.** *For each  $c > 0$ , we have, as  $x \rightarrow \infty$*

$$A(x, \varphi_M) = A(\varphi_M)x + O(x(\delta(x))^c)$$

where

$$A(\varphi_M) = \prod_{p \in \mathbf{P}} \left\{ 1 - p^{-s} + (1 - p^{-1}) \sum_{m=s}^{\infty} (\varphi_M(p))^{-1} \right\}$$

and  $\delta(x)$  is as given in (4.7).

REMARK 5.1. Theorem 5.1 affords us with another rich class of illustrations of Theorem 2.1. To this end, let  $t, K, r$  be integers such that  $t \geq 1$  and  $K > r \geq 2$ . We write

$$M_1 = M_1(r) = \{n \in \mathbb{Z}^+ | n \equiv r\},$$

$$M_2 = M_2(K, r) = \{n | n \text{ is congruent to one of } r, r+1, \dots, K-1 \pmod{K}\}$$

$$M_3 = M_3(t, r) = \{jr | j = 1, 2, \dots, t\},$$

$$M_4 = M_4(r) = \{jr | j \in \mathbb{Z}^+\}$$

and

$$M_5 = M_5(r) = \{r\}.$$

The elements of the sets  $Q_{M_1}$  through  $Q_{M_5}$  (usually denoted respectively by  $Q_r$ ,  $Q_{K,r}$ ,  $Q_{t,r}^*$ ,  $Q_r^*$  and  $Q_r^{s*}$ ) are known as  $r$ -free integers,  $(K, r)$ -integers ([22], [4]), unitarily  $(t, r)$ -integers, unitarily  $r$ -free integers ([3], [5]) and semi- $r$ -free integers ([23], [3], [5]). On specializing Theorem 5.1 with the set  $M$  chosen from among  $M_1, M_2, \dots, M_5$ , we obtain a number of illustrations of Theorem 2.1. It may be noted that the function  $\varphi_M$  in case  $M = M_1$  appears in [17] while in case  $M = M_2$  appears in [21] and [22].

## 6. An alternate form of Theorem 2.1

In this section, we give an equivalent form of Theorem 2.1 in that, in a way, it extends the sphere of applicability of Theorem 2.1.

If  $f, g$  are arithmetic functions, we write  $f * g$  to mean their Dirichlet product, that is, the arithmetic function defined by

$$(f * g)(n) = \sum_{d|n} f(d)g\left(\frac{n}{d}\right) \quad \text{for } n \in \mathbb{Z}^+.$$

Then we have the following result whose proof is simple.

THEOREM 6.1. Let  $f \in \mathcal{F}$  and  $g$  be a nonnegative multiplicative function satisfying

$$(6.1) \quad g(p^m) = 0 \quad \text{for } p \in P_1(f) \quad \text{and } m \in \mathbb{Z}^+;$$

$$(6.2) \quad \text{there exists an integer } l \text{ such that } g(p) = l \text{ for } p \in \mathbb{P} \setminus P_1(f).$$

Then the function  $f * g \in \mathcal{F}$  and as such the conclusion of Theorem 2.1 holds for  $f * g$  also.

REMARK 6.1. Suppose  $f \in \mathcal{F}$  with  $P_1(f) = \emptyset$ , the empty set. Then by Theorem 6.1, the conclusion of Theorem 2.1 is valid for the functions  $f * d_r, f * \theta$  etc. Here for  $r \in \mathbb{Z}^+$ ,  $d_r(n)$  denotes the Piltz divisor function defined to be the number of  $r$ -tuples

$(x_1, \dots, x_r)$  of positive integers satisfying  $x_1 x_2 \dots x_r = n$  and  $\theta(n)$  denotes the number of square-free divisors of  $n$ .

For the sake of simplicity, we considered in Theorem 2.1 functions  $f$  satisfying the condition (2.4). In a future communication, we extend Method C of Bateman [1] to a wide class of arithmetic functions  $f$  where  $f(p) = p^u + a_1 p^{u-1} + \dots + a_u$  for all primes  $p$  and  $a_1, \dots, a_u$  are rationals which do not depend on  $p$  and incidentally, we achieve results sharper than those contained in Theorems 4.1 and 5.1.

## REFERENCES

- [1] BATEMAN, P. T., The distribution of values of the Euler function, *Acta Arith.*, **21** (1972), 329—345. *MR* **46** # 1730.
- [2] COHEN, E., Arithmetical functions associated with the unitary divisors of an integer, *Math. Z.*, **74** (1960), 66—80. *MR* **22** # 3707.
- [3] COHEN, E., Some sets of integers related to the  $k$ -free integers, *Acta Sci. Math. (Szeged)* **22** (1961), 223—233. *MR* **24** # A1898.
- [4] COHEN, E., Arithmetical Notes XIII, A sequel to Note IV, *Elem. Math.*, **18** (1963), 8—11. *MR* **26** # 6104.
- [5] COHEN, E., Remark on a set of integers, *Acta Sci. Math. (Szeged)*, **25** (1964), 179—180. *MR* **30** # 66.
- [6] DIAMOND, H. G., Asymptotic distribution of Beurling's generalised integers, *Illinois J. Math.*, **14** (1970), 12—28. *MR* **40** # 5555.
- [7] DIAMOND, H. G., The distribution of values of Euler's phi function, *Analytic Number Theory (Proc. Sympos. Pure Math., Vol. 24 pp. 63—75. St. Louis Univ., St. Louis, Mo., 1972) Amer. Math. Soc., Providence, R. I., 1973. MR* **49** # 2604.
- [8] DICKSON, L. E., History of the theory of numbers, Vol. I, Chelsea Publishing Co. New York, 1966. *MR* **39** # 6807a.
- [9] DRESSLER, R. E., A density which counts multiplicity, *Pacific J. Math.*, **34** (1970), 371—378; *MR* **42** # 5940.
- [10] DRESSLER, R. E., An elementary proof of a theorem of Erdős on the sum of divisors function, *J. Number Theory*, **4** (1972), 532—536. *MR* **47** # 118.
- [11] ERDŐS, P., On the normal number of prime factors of  $p-1$  and some related problems concerning Euler's  $\phi$ -function, *Quart. J. Math., Oxford Ser. 6* (1935), 205—213. *Zbl* **12**. 149.
- [12] ERDŐS, P., Some remarks on Euler's  $\phi$ -function and some related problems, *Bull. Amer. Math. Soc.* **51** (1945), 540—544. *MR* **7**—49.
- [13] ESTERMANN, T., *Introduction to modern prime number theory*, Cambridge Tracts in Mathematics and Mathematical Physics, No. 41, Cambridge, at the University Press, 1952. *MR* **13**—915.
- [14] ISMAIL, M. and SUBBARAO, M. V., Unitary analogue of Carmichael's problem, *Indian J. Math.*, **18** (1976), 49—55. *MR* **82d**: 10011.
- [15] IVIČ, A., The distribution of values of some multiplicative functions, *Publ. Inst. Math. (Beograd, (N. S.))*, **22** (36) (1977), 87—94. *MR* **57** # 16229.
- [16] MALLIAVIN, P., Sur le reste de la loi asymptotique de répartition des nombres premiers généralisés de Beurling, *Acta Math.*, **106** (1961), 281—288. *MR* **26** # 87.
- [17] MCCARTHY, P. J., On a certain family of arithmetic functions, *Amer. Math. Monthly*, **65** (1958), 586—590. *MR* **20** # 6382.
- [18] NYMAN, B., A general prime number theorem, *Acta Math.*, **81** (1949), 299—307. *MR* **11**—332.
- [19] RIEGER, G. J., Einige verteilungsfragen mit  $K$ -leeren zahlen,  $r$ -zahlen und primzahlen, *J. Reine Angew. Math.*, **262/263** (1973), 189—193. *MR* **49** # 232.
- [20] SCHOENBERG, I. J., Über die asymptotische Verteilung reeller Zahlen mod 1, *Math. Z.*, **28** (1928), 171—199.
- [21] SUBBARAO, M. V., An arithmetic function and an associated probability theorem, *Nederl. Akad. Wetensch. Proc. Ser. A* **70** = *Indag. Math.* **29** (1967), 93—95. *MR* **34** # 7478.
- [22] SUBBARAO, M. V. and HARRIS, V. C., A new generalization of Ramanujan's sum, *J. London Math. Soc.*, **41** (1966), 595—604. *MR* **34** # 133.

- [23] SURYANARAYANA, D., Semi- $k$ -free integers, *Elem. Math.*, **26** (1971), 39—40. *MR* 47 # 147.
- [24] VIJAYARAGHAVAN, T., On the set of points  $\{\varphi(n)/n\}$ , *Jour. Indian Math. Soc.*, **12** (1920), 98—99.
- [25] WALFISZ, A., *Weylsche Exponentialsummen in der neueren Zahlentheorie*, Mathematische Forschungsberichte XV, VEB Deutscher Verlag der Wissenschaften, Berlin, 1963. *MR* 36 # 3737.

(Received December 20, 1983)

DEPARTMENT OF MATHEMATICS  
UNIVERSITY OF ALBERTA  
EDMONTON, ALBERTA  
CANADA

DEPARTMENT OF MATHEMATICS  
ANDHRA UNIVERSITY  
WALTAIR  
INDIA





# ON THE BOUNDEDNESS OF THE SOLUTIONS OF SECOND ORDER DIFFERENTIAL EQUATIONS

JERZY POPENDA

Applying Gronwall—Bihari—Bellman Lemma and their generalizations we study boundedness of the solutions of some ordinary differential equations. The method is similar to those used in [1], [4]. Furthermore we obtain conditions for the solutions to be of the square integrable with any wedge (see [4]).

LEMMA 1. Let  $u, f$  be nonnegative and continuous on  $I = \langle t_0, \infty \rangle$ ,  $C$ —any nonnegative constant. If

$$u(t) \leq C + \int_{t_0}^t f(s)u(s) ds \quad t \in I,$$

then

$$C + \int_{t_0}^t f(s)u(s) ds \leq C \exp \int_{t_0}^t f(s) ds \quad t \in I.$$

We denote by  $L_f^p(I)$  the class of all continuous functions  $u(t)$  on  $I$  such that

$$\int_{t_0}^{\infty} f(t)|u(t)|^p dt < \infty.$$

From Lemma 1 it follows that if  $f \in L_1^1(I)$ ,  $f$  and  $u$  are as in Lemma 1, then  $u(t)$  is bounded on  $I$  and  $u \in L_f^1(I)$ . It is obvious that if any function  $g(u)$  or any operator  $T(t, u)$  satisfy the inequality

$$u(t) \leq T(t, u)(g) \leq C + \int_{t_0}^t f(s)u(s) ds, \quad t \in I,$$

then

$$u(t) \leq T(t, u)(g) \leq C \exp \int_{t_0}^t f(s) ds, \quad t \in I.$$

Throughout the paper we assume that the considered equations have solutions of the class  $C^n(t_0, \infty)$ , where  $n$  denotes the order of the differential equation, and we shall study only such solutions. We assume furthermore that  $|x(t_0)| + \dots + |x^{(n-1)}(t_0)| > 0$ . We start with the equation of the form

$$(1) \quad x''(t) + a_1(t)x'(t) + a_0(t)x(t) + F(t, x(t), x'(t)) = e(t).$$

1980 *Mathematics Subject Classification*. Primary 34C10.

*Key words and phrases*. Differential equation, Gronwall lemma, boundedness, integrability.

**THEOREM 1.** Let  $a_1, a_0, a'_0, e$  be continuous on  $I$ ,  $a_0$  is positive,  $F(t, u, v)$  be continuous for  $t \in I$ ,  $u, v \in \mathbb{R}$ ;  $vF(t, u, v) \geq 0$  on  $I \times \mathbb{R}^2$ . If there exist positive continuous functions  $a$  and  $d$  on  $I$  such that

$$(2) \quad \int_{t_0}^{\infty} \frac{e^2(t)}{a^2(t)a_0(t)} dt < \infty, \quad \int_{t_0}^{\infty} r(t) dt < \infty,$$

where

$$(3) \quad r(t) = \max \left[ a^2(t) - 2a_1(t) - \frac{a'_0(t)}{a_0(t)} - d(t), d(t) \right],$$

then the solution  $x(t)$  is bounded,  $|x'(t)| = O(\sqrt{a_0(t)})$ ,  $x \in L^2_r(I)$ ,  $x' \in L^2_{r/a_0}(I)$ .

**PROOF.** Let  $D(t)$  be introduced by

$$(4) \quad D(t) = \frac{x'^2(t)}{a_0(t)} + x^2(t) + \int_{t_0}^t d(s)x^2(s) ds, \quad t \in I.$$

Then

$$(5) \quad D(t) = D(t_0) + \int_{t_0}^t \left[ \frac{2x'(s)x''(s)}{a_0(s)} - \frac{x'^2(s)a'_0(s)}{a_0^2(s)} + 2x(s)x'(s) + d(s)x^2(s) \right] ds, \quad t \in I.$$

By (1) we obtain

$$(6) \quad \begin{aligned} D(t) &= D(t_0) + \int_{t_0}^t \left[ \frac{2x'(s)e(s)}{a_0(s)} - \frac{2a_1(s)x'^2(s)}{a_0(s)} - \frac{a'_0(s)x'^2(s)}{a_0^2(s)} + \right. \\ &\quad \left. + d(s)x^2(s) - \frac{2x'(s)F(s, x(s), x'(s))}{a_0(s)} \right] ds \equiv \\ &\equiv D(t_0) + \int_{t_0}^t \frac{e^2(s)}{a^2(s)a_0(s)} ds + \int_{t_0}^t \left\{ \left[ a^2(s) - 2a_1(s) - \frac{a'_0(s)}{a_0(s)} \right] \frac{x'^2(s)}{a_0(s)} + d(s)x^2(s) \right\} ds \equiv \\ &\equiv D(t_0) + \int_{t_0}^{\infty} \frac{e^2(s)}{a^2(s)a_0(s)} ds + \\ &\quad + \int_{t_0}^t \left\{ \left[ a^2(s) - 2a_1(s) - \frac{a'_0(s)}{a_0(s)} - d(s) \right] \frac{x'^2(s)}{a_0(s)} + d(s)D(s) \right\} ds, \quad t \in I. \end{aligned}$$

From (2), (3) and (4) we have

$$D(t) \leq M + \int_{t_0}^t 2r(s)D(s) ds, \quad t \in I,$$

for any positive constant  $M$ . Now applying Lemma 1 and (4) we obtain boundedness of the solution and its derivative. It suffices to observe that the rest of the theorem follows from the remark to the Lemma 1.

REMARK 1. Observe that Theorem 1 remains true for more general equations, in which the function  $F(t, u, v)$  is replaced by  $F(t, u, v, \dots)$ . This means that in the equation (1) we may have the term of the form  $F(t, x(t), x'(t), \dots)$ . Let the function  $F$  be defined and continuous on  $I \times \mathbb{R}^m$  and satisfies the inequality  $vF(t, u, v, \dots) \geq 0$ .

DEFINITION. We say that the function  $F$  is of the class  $B_{A,C}^{1,n}$  if

$$(i) \quad |F(t, x_1, \dots, x_n)| \leq C(t, |x_1|, \dots, |x_n|) \quad \text{for } t \geq t_0 \geq 0, \\ x_1, \dots, x_n \in \mathbb{R},$$

$$(ii) \quad C(t, x_1, \dots, x_n) \leq C(t, \bar{x}_1, \dots, \bar{x}_n) \quad \text{for } t \geq t_0 \geq 0, \\ \bar{x}_i \geq x_i \geq 0, \quad i = 1, \dots, n,$$

$$(iii) \quad C(t, a(t)x_1, \dots, a(t)x_n) \leq A(a(t))C(t, x_1, \dots, x_n)$$

for arbitrary continuous, positive function  $a(t), t \in I$ , where  $A$  is positive continuous,  $x_i \in \mathbb{R}_+, i = 1, \dots, n$ .

LEMMA 2. Let  $u, v$  be nonnegative and continuous on  $I, C$ —any positive constant,  $f$  continuous and positive on  $(C, \infty)$ , nondecreasing. If

$$u(t) \leq C + \int_{t_0}^t v(s)f(u(s))ds, \quad t \in I,$$

then

$$C + \int_{t_0}^t v(s)f(u(s))ds \leq F^{-1}\left(F(C) + \int_{t_0}^t v(s)ds\right), \quad t \in I,$$

where  $F^{-1}$  denotes the inverse function of  $F(t) = \int_{\varepsilon}^t \frac{1}{f(s)}ds, t \geq \varepsilon, \varepsilon$  is any fixed positive constant.

We can make here the same remarks as to Lemma 1.

THEOREM 2. Let  $a_1, a_0, a'_0, e$  be continuous on  $I, a_0$  is positive,  $F(t, u, v)$  is continuous for  $t \in I, u, v \in \mathbb{R}, F^2 \in B_{A,C}^{1,2}$  and  $\lim_{u \rightarrow \infty} G(u) = \infty$  where  $G(u) = \int_{\varepsilon}^u \frac{dv}{v + A(\sqrt{v})}$ . If there exist continuous, positive functions  $a, d$  on  $I$  such that

$$(7) \quad \int_{t_0}^{\infty} \frac{e^2(t)}{a^2(t)a_0(t)}dt < \infty, \quad \int_{t_0}^{\infty} r(t)dt < \infty,$$

where

$$(8) \quad r(t) = \max \left[ a^2(t) + d^2(t) - 2a_1(t) - \frac{a'_0(t)}{a_0(t)}, \frac{C(t, 1, \sqrt{a_0(t)})}{d^2(t)a_0(t)} \right]$$

then every solution  $x(t)$  of the class  $C^2(t_0, \infty)$  to the differential equation (1) is bounded on  $I$ , and

$$|x'(t)| = O(\sqrt{a_0(t)}), \quad \text{as } t \rightarrow \infty, \quad x \in L^2_r(I), x' \in L^2_{r/a_0}(I).$$

PROOF. Let  $D(t)$  be defined by

$$D(t) = \frac{x'^2(t)}{a_0(t)} + x^2(t), \quad t \in I.$$

Then analogously as in Theorem 1 we have

$$\begin{aligned} D(t) &\equiv D(t_0) + \int_{t_0}^t \frac{e^2(s)}{a^2(s)a_0(s)} ds + \int_{t_0}^t \left[ a^2(s) - 2a_1(s) - \frac{a'_0(s)}{a_0(s)} \right] \frac{x'^2(s)}{a_0(s)} + \\ &\quad + 2 \frac{|x'(s)|}{\sqrt{a_0(s)}} \frac{|F(s, x(s), x'(s))|}{\sqrt{a_0(s)}} ds \equiv D(t_0) + \int_{t_0}^\infty \frac{e^2(s)}{a^2(s)a_0(s)} ds + \\ &\quad + \int_{t_0}^t \left\{ \left[ a^2(s) - 2a_1(s) - \frac{a'_0(s)}{a_0(s)} + d^2(s) \right] \frac{x'^2(s)}{a_0(s)} + \frac{C(s, |x(s)|, |x'(s)|)}{d^2(s)a_0(s)} \right\} ds \equiv \\ &\equiv D(t_0) + \int_{t_0}^\infty \frac{e^2(s)}{a^2(s)a_0(s)} ds + \int_{t_0}^t \left[ r(s) \frac{x'^2(s)}{a_0(s)} + A(\sqrt{D(s)}) \frac{C(s, 1, \sqrt{a_0(s)})}{d^2(s)a_0(s)} \right] ds, \quad t \in I. \end{aligned}$$

Hence

$$\begin{aligned} D(t) &\equiv D(t_0) + \int_{t_0}^t r(s) x^2(s) ds \equiv \\ &\equiv D(t_0) + \int_{t_0}^\infty \frac{e^2(s)}{a^2(s)a_0(s)} ds + \int_{t_0}^t r(s) [D(s) + A(\sqrt{D(s)})] ds, \quad t \in I. \end{aligned}$$

Now, applying Lemma 2 we obtain the desired estimations announced in the thesis.

REMARK 2. Theorem 2 holds for the equations of more general type. In these types the function  $F$  in equation (1) is replaced by the other of more variables. Conditions (7), (8) ought to be changed and used other generalizations of Gronwall lemma (see [2], [3]). This generalization can take the following form.

Let  $a_1, a_0, a'_0, e$  be continuous, and  $a_0$  be positive on  $I$ ,  $F(t, u, v, w)$  be continuous for  $t \in I, u, v, w \in \mathbb{R}$ ,  $F \in B_{\lambda, \epsilon}^{1,3}$  and  $\lim_{u \rightarrow \infty} G(u) = \lim_{u \rightarrow \infty} \int_{\epsilon}^u \frac{dt}{\sqrt{t} + \sqrt{t} A(\sqrt{t}) + t} = \infty$ . Let the function  $b$  be continuous on  $I$  satisfying  $\lim_{t \rightarrow \infty} b(t) = \infty$ ,  $b(t) \leq t$ . If furthermore  $\int_{t_0}^\infty r(t) dt < \infty$ , where

$$r(t) = \max \left[ \frac{2|e(t)|}{\sqrt{a_0(t)}}, \frac{2C(t, 1, \sqrt{a_0(t)}, \sqrt{a_0(b(t))})}{\sqrt{a_0(t)}}, -2a_1(t) - \frac{a'_0(t)}{a_0(t)} \right].$$

Then the solution of the equation

$$x''(t) + a_1(t)x'(t) + a_0(t)x(t) + F(t, x(t), x'(t), x'(b(t))) = e(t)$$

is bounded,  $|x'(t)| = O(\sqrt{a_0(t)})$ ,  $x \in L_r^2(I)$ ,  $x' \in L_{r/a_0}^2(I)$ . The present method can be

applied to study other properties of the solutions. We give one example formulated in the following theorem similar to Theorem 1.

**THEOREM 3.** Let  $a_1, a_0, a'_0, e$  be continuous functions on  $I, a_0 > 0$  there. Let  $F(t, u, v)$  be continuous for  $t \in I, u, v \in \mathbb{R}$  and  $vF(t, u, v) \geq 0$  on  $I \times \mathbb{R}^2$ . If there exist positive, continuous functions  $a$  and  $d$  on  $I$ , such that  $d$  is continuously differentiable,  $d' \geq 0, \lim_{t \rightarrow \infty} d(t) = \infty$ , and

$$\int_{t_0}^{\infty} \frac{e^2(t)d(t)}{a^2(t)a_0(t)} dt < \infty, \quad \int_{t_0}^{\infty} \left[ a^2(t) - 2a_1(t) - \frac{d(t)}{a_0(t)} \left( \frac{a_0(t)}{d(t)} \right)' \right]^+ dt < \infty.$$

Then the solution  $x(t)$  of the equation (1) tends to zero as  $t \rightarrow \infty$ . (Here  $f^+(t) = \max[f(t), 0]$ ).

The considered method and theorems can be extended and applied for certain differential equations of higher order.

#### REFERENCES

- [1] BEHZAD, M. and MEHRI, B., On boundedness property of the solution of certain nonlinear differential equations, *Studia Sci. Math. Hungar.* 6 (1971), 163—166. MR 48 #6566.
- [2] BIHARI, I., Researches of the boundedness and stability of the solutions of non-linear differential equations, *Acta Math. Acad. Sci. Hungar.* 8 (1957), 261—278. MR 20 #1031.
- [3] BOBROWSKI, D., POPENDA, J. and WERBOWSKI, J., On the systems integral inequalities with delay of Gronwall—Bellman type, *Fasc. Math.* 10 (1978), 97—104. MR 58 #17010.
- [4] WYRWIŃSKA, A., O rozwiązaniach pewnych równań różniczkowych w klasie  $L_p$ , *Doctoral dissertation*, Politechnika Poznańska, Poznań, 1977 (Polish).

(Received December 16, 1980)

INSTYTUT MATEMATYCZNY  
POLITECHNIKA POZNAŃSKA  
PIOTROWO 3A  
PL-60-965 POZNAŃ  
POLAND





# ПАНЦИКЛИЧНОСТЬ ОРГРАФОВ ПРИ УСЛОВИИ МЕЙНИЛА

С. Х. ДАРБИНЯН

## Abstract

Let  $G$  be a strongly connected digraph with  $n$  vertices satisfying the condition that the sum of degrees for any two distinct nonadjacent vertices is at least  $2n-1$ . Then either

- 1)  $G$  is pancyclic;
- or 2)  $n$  is even and  $G \in \{\bar{K}_{n/2, n/2}, \bar{K}_{n/2, n/2} \setminus \{e\}\}$ , where  $e \in E(\bar{K}_{n/2, n/2})$ .
- or 3)  $G \in \Phi_n^m$ , where  $n-1 \equiv m \pmod{(n+1)/2}$  and  $\Phi_n^m$  is the family of digraphs  $G'$ , which satisfies the following conditions:

- a)  $V(G') = \{x_1, x_2, \dots, x_n\}$ ;
- b) for each  $i \in [1, n-m+1]$  the vertices  $x_i$  and  $x_{i+m-1}$  are nonadjacent;
- c)  $x_1, x_n \in E(G')$  and  $x_{i+1}x_i \in E(G')$  for each  $i \in [1, n-1]$ ;
- d) if  $2 \leq i+1 < j \leq n$ , then  $x_jx_i \in E(G')$ ;
- e) the sum of degrees for any two distinct nonadjacent vertices is at least  $2n-1$ .

Будем рассматривать конечные орграфы (графы) без петель и кратных дуг (ребер). Все понятия и обозначения, не определяемые здесь, можно найти в книге [1]. Орграф (граф) с  $n$  вершинами ( $n \geq 3$ ) называется панциклическим, если он содержит контур (цикл) любой длины  $k$ , при  $3 \leq k \leq n$ .

Бонди [3] была предложена следующая

Метагипотеза. Почти все «нетривиальные» достаточные условия гамильтоновости графа являются достаточными условиями для того, чтобы граф был панциклическим (быть может кроме «простого» класса графов).

Будем говорить, что  $n$ -вершинный ( $n \geq 2$ ) орграф удовлетворяет условию  $(N_i)$ ,  $0 \leq i \leq 2$ , если сумма степеней любых двух несмежных различных вершин не меньше  $2n-1+i$ .

Мейнил [4] доказал, что если сильно связный орграф удовлетворяет условию  $(N_0)$ , то он гамильтонов. Возникает вопрос: удовлетворяет ли достаточное условие гамильтоновости орграфов Мейнила метагипотезе Бонди? В настоящей работе характеризуется класс тех сильно связных орграфов, которые удовлетворяют условию  $(N_0)$  и не являются панциклическими.

Через  $V(G)$  обозначим множество вершин орграфа  $G$ , а через  $E(G)$  — множество его дуг. Дуга исходящая из вершины  $u$  в вершину  $v$ , обозначим через  $uv$ . Пусть  $A, B \subseteq V(G)$  и  $x \in V(G)$ . Введем обозначения:

$$E(A \rightarrow B) = \{zy \in E(G) / z \in A, y \in B\};$$

$$E(A, B) = E(A \rightarrow B) \cup E(B \rightarrow A);$$

$$d(x) = |E(\{x\}, V(G))|; \quad d(x, A) = |E(\{x\}, A)|.$$

1980 Mathematics Subject Classification. Primary 05C20; Secondary 05C38.

Key words and phrases. Digraph, cycle, Hamilton cycle, strongly connected digraph, pancyclic digraph.

Запись  $A \rightarrow B$  означает, что если  $y \in A$  и  $z \in B$ , то  $yz \in E(G)$ . Если  $H \subseteq V(G)$  и  $A \rightarrow B$ ,  $B \rightarrow H$ , то будем писать  $A \rightarrow B \rightarrow H$ . Если же  $A = \{x\}$ , то вместо  $\{x\}$  будем писать  $x$ . Подграф порожденный множеством вершин  $A$  обозначим через  $\langle A \rangle$ .

Путь, состоящий из вершин  $x_1, x_2, \dots, x_k$  и дуг  $x_i x_{i+1}$ ,  $1 \leq i \leq k-1$  обозначим через  $(x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_k)$ , а контур, полученный из этого пути после добавления дуги  $x_k x_1$  — через  $(x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_k \rightarrow x_1)$ .

Путь  $(x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_k)$ , называем  $(A, B)$ -путем, если  $x_1 \in A$ ,  $x_k \in B$  и для всех  $i$  ( $2 \leq i \leq k-1$ ) имеет место  $x_i \notin A \cup B$ . Будем говорить, что дуга  $xu \in E(G)$  входит в  $A$  (исходит из  $A$ ), если  $x \notin A$  и  $y \in A$  (соответственно  $x \in A$  и  $y \notin A$ ).

Для пути  $P: (x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_k)$  введем обозначения:

$$P[x_i, x_j]: (x_i \rightarrow x_{i+1} \rightarrow \dots \rightarrow x_j), \quad 1 \leq i \leq j \leq k;$$

$$P(x_i, x_j): (x_{i+1} \rightarrow x_{i+2} \rightarrow \dots \rightarrow x_j), \quad 1 \leq i < j \leq k;$$

$$P[x_i, x_j): (x_i \rightarrow x_{i+1} \rightarrow \dots \rightarrow x_{j-1}), \quad 1 \leq i < j \leq k;$$

(иногда вершину орграфа  $G$  будем считать путем длины 0 и контуром длины 1). Аналогичными обозначениями будем пользоваться для контуров.  $\bar{K}_{n,n}$  — симметрический двудольный орграф полученный из  $K_{n,n}$ .

Пусть  $P: (x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_m)$  и  $Q: (y_1 \rightarrow y_2 \rightarrow \dots \rightarrow y_k)$  некоторые пути в орграфе  $G$ . Если  $V(P) \cap V(Q) = \emptyset$  и  $x_m y_1 \in E(G)$ , то обозначим через  $(P \rightarrow Q)$  путь  $(x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_m \rightarrow y_1 \rightarrow y_2 \rightarrow \dots \rightarrow y_k)$ , а если же вершина  $x_m$  совпадает с вершиной  $y_1$ , а вершина  $y_k$  с вершиной  $x_1$  и внутренние вершины этих путей не совпадают друг с другом, тогда  $(P \rightarrow Q)$  означает контур  $(x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_m \rightarrow y_2 \rightarrow y_3 \rightarrow \dots \rightarrow y_k)$ .

Пусть орграф  $G$  содержит гамильтоновы контур  $H: (x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_n \rightarrow x_1)$ . Через  $H[x_i]$  (соответственно,  $H(x_i)$ ) обозначим гамильтоновы путь  $(x_i \rightarrow x_{i+1} \rightarrow \dots \rightarrow x_{i-1})$  (соответственно,  $(x_{i+1} \rightarrow x_{i+2} \rightarrow \dots \rightarrow x_{i-1} \rightarrow x_i)$ ), а через  $H^*$  — гамильтоновы путь графа  $G$ , полученный из контура  $H$  после удаления одной дуги.

Как обычно  $C_k$  означает контур длины  $k$ , а  $[m, n]$  множество целых чисел, не больших  $n$  и не меньших  $m$ . Если  $C_k: (x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_k \rightarrow x_1)$ , то всюду индексы вершин контура  $C_k$  берутся по mod  $(k)$ .

В дальнейшем нам понадобятся следующие утверждения:

**Предложение 1 ([5]).** Пусть  $n$ -вершинный орграф  $G$  содержит контур  $C_k: (x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_k \rightarrow x_1)$ , где  $2 \leq k \leq n-1$ . Если для некоторой вершины  $x \in V(G) \setminus V(C_k)$  имеет место  $d(x, V(C_k)) \equiv k+1$ , то для любого  $m \in [2, k+1]$  орграф  $G$  содержит контур  $C_m$  такой, что  $V(C_m) \subseteq \{x\} \cup V(C_k)$ .

**Предложение 2 ([6]).** Пусть  $n$ -вершинный орграф  $G$  содержит путь  $P: (x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_k)$ , где  $2 \leq k \leq n-1$ . Если для некоторой вершины  $x \in V(G) \setminus V(P)$  имеет место

$$d(x, V(P)) \equiv \begin{cases} k+1, & \text{если } x_k x \notin E(G) \text{ или } x x_1 \notin E(G), \\ k+2, & \text{если } x_k x \in E(G), \quad x x_1 \in E(G) \end{cases}$$

то в  $G$  существует  $(x_1, x_k)$ -путь, проходящий через все вершины  $x_1, x_2, \dots, x_k$  и  $x$  (т.е. путь  $P$  можно расширить с помощью вершины  $x$ ).

Всюду в работе компоненты  $G_1, G_2, \dots, G_s$  орграфа  $G$  будут пронумерованы так, что нет дуг, исходящий из компоненты с большим номером и входящий в компоненту с меньшим номером;  $G_1$  назовем начальной компонентой, а  $G_s$  — конечной.

**Предложение 3 ([7]).** Пусть  $T$  —  $n$ -вершинный ( $n \geq 2$ ) не сильно связный турнир и  $T^1, T^2, \dots, T^s$  его компоненты. Тогда для любых вершины  $x \in V(T^1)$ ,  $y \in V(T^s)$  и для любого числа  $l$  ( $1 \leq l \leq n-1$ ) существует  $(x, y)$ -путь длины  $l$ . Более того, если  $2 \leq l \leq n-1$ , то  $(x, y)$ -путь можно выбрать так, чтобы он проходил через данную вершину из  $T^i$ , где  $2 \leq i \leq s-1$ .

**Лемма 1.** Пусть в  $n$ -вершинном ( $n \geq 1$ ) орграфе  $G$  всякие две различные вершины смежны между собой. Тогда

а)  $G$  содержит остовный подграф, являющийся турниром.

б) если  $n \geq 3$ , то  $G$  является панциклическим тогда и только тогда, когда  $G$  — сильно связный.

**Доказательство.** Утверждение а) и необходимость утверждения б) очевидны. А достаточность утверждения б) доказывается точно также, как теорема 16 II в книге [1].

Для любого  $n \geq 5$  и  $m$  ( $n-1 \geq m > (n+1)/2$ ) обозначим через  $\Phi_n^m$  множество всевозможных  $n$ -вершинных орграфов  $G$ , удовлетворяющие условию  $(N_0)$ , для которых имеют место

I.  $V(G) = \{x_1, x_2, \dots, x_n\}$ ;

II.  $x_1 x_n \in E(G)$  и  $x_{i+1} x_i \in E(G)$ , при всех  $i \in [1, n-1]$ ;

III.  $E(x_i, x_{i+m-1}) = \emptyset$ , при всех  $i \in [1, n-m+1]$ ;

IV. если  $2 \leq i+1 < j \leq n$ , то  $x_j x_i \notin E(G)$ .

Пусть  $G$  есть  $n$ -вершинный ( $n \geq 3$ ) орграф. Обозначим через  $C(G)$  множество таких контуров  $C_m \subset G$  ( $2 \leq m \leq n-1$ ) для которых степень каждой вершины, не принадлежащий контуру  $C_m$ , не меньше  $n$ .

**Лемма 2.** Пусть  $n$ -вершинный ( $n \geq 3$ ) сильно связный орграф  $G$  удовлетворяет условию  $(N_0)$ . Тогда справедливо хотя бы одно из следующих утверждений: а)  $C(G) \neq \emptyset$ ; б)  $G$  — панциклический; в) для некоторого  $m$  ( $n-1 \geq m > (n+1)/2$ )  $G \in \Phi_n^m$ .

**Доказательство.** При  $n=3$  утверждение леммы очевидно. Пусть  $n=4$  и

$$G_0 = \langle \{x \in V(G) / d(x) \leq n-1\} \rangle; \quad Q = \langle V(G) \setminus V(G_0) \rangle$$

и  $G_1, G_2, \dots, G_s$  компоненты орграфа  $G_0$ .

Пусть  $p = |V(G_0)|$  и  $q = |V(Q)|$ .

Пусть  $p \leq 1$ . Если  $p=0$ , то  $G$  содержит контур  $C_2 \in C(G)$ . Поэтому будем предполагать, что  $p=1$  и пусть  $\{x\} = V(G_0)$ . Из сильно связности  $G$  вытекает существование таких вершин  $y, z \in V(Q)$ , что  $xy, zx \in E(G)$ . Так как вершина  $z$  достижима из вершины  $y$ , то в  $G$  существует  $(y, z)$ -путь  $P: (y = u_1 \rightarrow u_2 \rightarrow \dots \rightarrow u_k = z)$ . Предположим, что  $C(G) = \emptyset$ . Тогда нетрудно убедиться, что  $y \neq z$ ;  $V(P) = V(Q)$  и  $u_i u_j \in E(G)$ , при  $1 \leq i < j \leq k$ , тогда и только тогда, когда  $j = i+1$ . Следовательно, поскольку  $d(u_1) \geq n$ , то  $\{u_2, u_3, \dots, u_k\} \rightarrow u_1$  и, значит,  $G$  является панциклическим. Поэтому в дальнейшем будем предполагать, что  $p \geq 2$ .

Из условия  $(N_0)$  вытекает, что любые две различные вершины множества  $V(G_0)$  смежны между собой. Поэтому для всех  $i$  и  $j$  ( $1 \leq i < j \leq s$ ) имеет место

$$(1) \quad V(G_i) \rightarrow V(G_j)$$

и по лемме 1 любой порожденный подграф  $\langle B \rangle$ , где  $B \subseteq V(G_0)$  содержит остовный подграф, являющийся турниром. Причем, если  $\langle B \rangle$  является сильно связным, то при  $|B| \geq 3$ , подграф  $\langle B \rangle$  является панциклическим. Следовательно, по лемме 1 каждая компонента  $G_i$  ( $1 \leq i \leq s$ ) является гамильтоновой (здесь и ниже одновершинные графы считаются гамильтоновыми). Отсюда получаем, что если  $G_0$  является сильно связным, то по лемме 1 имеем: если  $V(Q) = \emptyset$ , то  $G$  является панциклическим, а если же  $V(Q) \neq \emptyset$ , то  $C(G) \neq \emptyset$ . Поэтому будем предполагать, что  $G_0$  не является сильно связным, т.е.  $s \geq 2$ . Следовательно, поскольку  $G$  — сильно связный, то  $V(Q) \neq \emptyset$  и существует  $(V(G_s), V(G_1))$ -путь. Выберем в каждом  $G_i$ ,  $1 \leq i \leq s$ , некоторый гамильтоновый контур  $H_i$  (если  $G_i$  одновершинный, то считаем, что  $H_i = x$ , где  $\{x\} = V(G_i)$ ) и фиксируем эти контуры. Обозначим через  $R$  множество таких всевозможных  $(V(G_s), V(G_1))$ -путей  $P$ , для которых имеют место следующие условия:

1) не более одной дуги пути  $P$  входит в  $V(G_i)$  и не более одной дуги пути  $P$  исходит из  $V(G_i)$ ,  $2 \leq i \leq s-1$ .

2) если  $V(P) \cap V(G_i) \neq \emptyset$ , то множество дуг  $E(P) \cap E(G_i)$  составляют подпуть для  $H_i$ , т.е. в  $V(G_i)$  путь  $P$  проходит по направлению контура  $H_i$ .

Пусть  $(V(G_s), V(G_1))$ -путь  $P: (u_1 \rightarrow u_2 \rightarrow \dots \rightarrow u_i)$  такой, что количество тех компонент подграфа  $G_0$ , для которых не выполняется условия 1) или 2) минимальный и равно  $\alpha$ . Покажем, что  $\alpha = 0$ . Предположим, что  $\alpha \geq 1$ . Пусть  $S = V(P) \cap V(G_i) \neq \emptyset$ , где  $2 \leq i \leq s-1$ , и для  $G_i$  не выполняется условия 1) или 2). В пути  $P$  подпуть  $P[u_{j_1}, u_{j_2}]$  заменим на подпуть  $H_i[u_{j_1}, u_{j_2}]$ , где  $u_{j_1}, u_{j_2} \in S$  и  $u_{j_1}$  (соответственно  $u_{j_2}$ ) та вершина, которая имеет минимальный (соответственно максимальный) индекс. В результате получаем  $(V(G_s), V(G_1))$ -путь  $P_1$  с количеством компонент подграфа  $G_0$ , для которых не выполняется условия 1) или 2), меньшим  $\alpha$ , а это противоречит минимальности  $\alpha$ . Итак,  $\alpha = 0$ .

Равенство  $\alpha = 0$  означает, что  $R \neq \emptyset$ .

Так как для всех  $(z, y)$ -путей  $P \in R$  подграф  $\langle V(G_0) \setminus V(P(y, z)) \rangle$  содержит остовный подграф, являющийся не сильно связным турниром, в котором вершина  $y$  принадлежит начальной компоненте, а вершина  $z$  — конечной компоненте, то по предложению 3 он содержит гамильтоновый  $(y, z)$ -путь  $F$ . Отсюда следует, что  $(P \rightarrow F)$  является гамильтоновым контуром для подграфа  $\langle V(G_0) \cup V(P) \rangle$ . Следовательно, если существует путь  $P \in R$  не проходящий через все вершины множества  $V(Q)$ , то контур  $(P \rightarrow F)$  принадлежит множеству  $C(G)$ . Значит в этом случае  $C(G) \neq \emptyset$ . Поэтому в дальнейшем будем предполагать, что всякий путь из  $R$  проходит через все вершины  $V(Q)$ , т.е. если  $P \in R$ , то

$$(2) \quad V(Q) \setminus V(P) = \emptyset.$$

Пусть  $P \in R$  и  $P^1, P^2, \dots, P^v$  такие максимальные по длине подпути для  $P$ , что для любого  $i \in [1, v]$  имеет место  $V(P^i) \subseteq V(Q)$  и подпути  $P^1, P^2, \dots, P^v$  пронумерованы таким образом, что при любых  $1 \leq i < j \leq v$ , по направлению пути  $P$  вершины множества  $V(P^i)$  встречаются после вершины множества  $V(P^j)$ . Подпути  $P^1, P^2, \dots, P^v$  назовем  $P(Q)$ -подпутями. Пусть для определен-

ности  $(z, y)$ -путь  $P \in R$  такой, что количество  $P(Q)$ -подпутей наименьшее и равно  $v$ .

Для любого  $k \in [1, v]$  пусть

$$m_k = |V(P^k)|; \quad f(k) = \sum_{i=1}^k m_i; \quad f(0) = 0;$$

и пусть для определенности  $P^k: (x_{f(k)} \rightarrow x_{f(k)-1} \rightarrow \dots \rightarrow x_{f(k-1)+1})$ . Тогда, в частности, имеем  $V(Q) = \{x_1, x_2, \dots, x_{f(v)}\}$  и  $q = f(v)$ .

В ходе доказательства леммы 2 будут доказаны ряд свойств в виде утверждений.

Утверждение 1. а) Если  $1 \leq i < j \leq v$ , то

$$L_1 \triangleq E(V(P^j) \rightarrow V(P^i)) = \emptyset.$$

б) Если  $1 \leq t < k-1 \leq q-1$ , то  $x_k x_t \notin E(G)$ .

Доказательство. а) Допустим, что  $L_1 \neq \emptyset$  и пусть  $uv \in L_1$ . Тогда путь  $P_1: (P[z, u] \rightarrow P[v, y])$  принадлежит множеству  $R$  и количество  $P_1(Q)$ -подпутей меньше  $v$ , а это противоречит минимальности  $v$ .

б) Предположим, что  $x_k x_t \in E(G)$ . Тогда путь  $P_2: (P[z, x_k] \rightarrow P[x_t, y])$  принадлежит множеству  $R$  и  $x_{t+1} \in V(Q) \setminus V(P_2)$ , т. е.  $V(Q) \setminus V(P_2) \neq \emptyset$ , а это противоречит соотношению (2).

Утверждение 2. Для любого  $j \in [1, q]$  имеет место

$$d(x_j, V(Q)) \leq \begin{cases} q-1, & \text{если } \{x_j\} = V(P^i), \quad 1 \leq i \leq v; \\ q, & \text{если вершина } x_j \text{ является началом или концом подпути } P^j, \quad 1 \leq j \leq v; \\ q+1, & \text{в остальных случаях.} \end{cases}$$

Доказательство следует из утверждения 1.

Утверждение 3.

$$L_2 \triangleq E(V(Q) \setminus \{x_1\} \rightarrow V(G_1)) = \emptyset;$$

$$L_3 \triangleq E(V(G_s) \rightarrow V(Q) \setminus \{x_q\}) = \emptyset.$$

Действительно, если  $L_2 \neq \emptyset$  и  $u_1 v_1 \in L_2$ , то путь  $P_1: (P[z, u_1] \rightarrow v_1) \in R$  и  $x_1 \notin V(P_1)$ , если же  $L_3 \neq \emptyset$  и  $u_2 v_2 \in L_3$ , то путь  $P_2: (u_2 \rightarrow P[v_2, y]) \in R$  и  $x_q \notin V(P_2)$ , а это противоречит (2), что и доказывает утверждение 3.

Из максимальности подпутей  $P^1, P^2, \dots, P^v$  вытекает, что если  $v \geq 2$ , то для любого  $i \in [1, v-1]$  существует такие числа  $i_1, t_i \in [2, s-1]$ , что

$$E(x_{f(i)+1} \rightarrow V(G_{i_1})) \neq \emptyset; \quad E(V(G_{t_i}) \rightarrow x_{f(i)}) \neq \emptyset.$$

Выберем минимальное  $i_1$  и максимальное  $t_i$  удовлетворяющие этим требованиям. Так как  $E(V(G_{j_2}) \rightarrow V(G_{j_1})) = \emptyset$  ( $1 \leq j_1 < j_2 \leq s$ ), то для всех  $i \in [1, v-1]$  имеем  $i_1 \leq t_i$ .



Докажем, что если  $j \in [1, v-2]$ , то  $t_j < t_{j+1}$ . Предположим обратное. Пусть

$$u \in V(G_{i_{j+1}}); \quad v \in V(G_{t_j}); \quad x_{f(j+1)+1}u, vx_{f(j)} \in E(G)$$

и

$$M_1 = V(P) \cap V(G_{i_{j+1}}); \quad M_2 = V(P) \cap V(G_{t_j}).$$

а)  $M_1 = M_2 = \emptyset$ . Тогда путь  $P_1: (P[z, x_{f(j+1)+1}] \rightarrow K_1 \rightarrow P[x_{f(j)}, y]) \in R$  и  $V(P^{j+1}) \not\subseteq V(P_1)$ , где  $K_1: (u \rightarrow v)$  или  $K_1: H_{i_{j+1}}[u, v]$  для  $t_j > i_{j+1}$  и  $t_j = i_{j+1}$ , соответственно. А это противоречит соотношению (2).

б)  $M_1 \neq \emptyset$  и  $M_2 = \emptyset$ . Пусть вершины множества  $M_1$  на пути  $P$  встречаются после вершин  $V(P^i)$  и раньше вершин  $V(P^{i-1})$ . Если  $i \leq j+1$ , то путь  $P_2: (P[z, x_{f(j+1)+1}] \rightarrow K_2 \rightarrow P[x_{f(i-1)}, y]) \in R$  и  $V(P^{j+1}) \not\subseteq V(P_2)$ , где  $v_1 x_{f(i-1)} \in E(P)$  и  $K_2: (u \rightarrow v_1)$  или  $K_2: H_{i_{j+1}}[u, v_1]$  для  $v_1 \notin V(G_{i_{j+1}})$  и  $v_1 \in V(G_{i_{j+1}})$ , соответственно. Если же  $i \geq j+2$ , то путь  $P_3: (P[z, x_{f(i-1)+1}] \rightarrow K_3 \rightarrow P[x_{f(j)}, y]) \in R$  и  $V(P^{j+1}) \not\subseteq V(P_3)$ , где  $x_{f(i-1)+1}u_1 \in E(P)$  и  $K_3: (u_1 \rightarrow v)$ .

Итак, вновь пришли к противоречию.

Случаи в)  $M_1 = \emptyset$  и  $M_2 \neq \emptyset$  и г)  $M_1 \neq \emptyset$  и  $M_2 \neq \emptyset$  рассматриваются аналогичным образом.

Из  $i_j \leq t_j$  и  $t_j < i_{j+1}$  следует, что  $i_1 \leq t_1 < i_2 \leq t_2 < \dots < i_{v-1} \leq t_{v-1}$ .

Для любого  $i \in [1, v-1]$  определим пути  $S_i$  следующим образом: если  $i < t_i$ , то  $S_i: (H_i[u] \rightarrow H_{i+1}^* \rightarrow \dots \rightarrow H_{t_i-1}^* \rightarrow H_{t_i}(v))$ , а если же  $i = t_i$ , то  $S_i: H_i[u, v]$  и вершины  $u$  и  $v$  выбраны таким образом, чтобы путь  $S_i$  имел возможно большую длину, где  $u \in V(G_{i_i})$ ;  $v \in V(G_{t_i})$  и  $x_{f(i)+1}u, vx_{f(i)} \in E(G)$ .

Пусть  $S_0 = H_1[y]$ ;  $S_v = H_v(z)$  и  $s_j = |V(S_j)|$ ,  $j \in [0, v]$ .

Из всего этого непосредственно вытекает, что путь  $(z \rightarrow P^v \rightarrow S_{v-1} \rightarrow \dots \rightarrow P^i \rightarrow S_{i-1} \rightarrow P^{i-1} \rightarrow \dots \rightarrow P^1 \rightarrow y)$  принадлежит множеству  $R$ . Полученный путь вновь обозначим через  $P$ . Очевидно, что количество  $P(Q)$ -подпутей равно  $v$ .

Из определений подпутей  $S_j$  ( $0 \leq j \leq v$ ) и соотношения (1) следует, что можно выбрать такой гамильтоновы путь  $H: (y = y_1 \rightarrow y_2 \rightarrow \dots \rightarrow y_p = z)$  в  $G_0$ , который проходить по направлению фиксированных гамильтоновых контуров  $H_i$  ( $1 \leq i \leq s$ ); подпути  $S_j$  являются подпутями для  $H$  и если  $V(S_j) \subseteq V(G_l)$  ( $2 \leq l \leq s-1$ ), то путь  $H$  из  $V(G_{l-1})$  сразу входит в начальную вершину подпути  $S_j$ .

Поскольку подграф  $\langle V(G_0) \setminus V(P(y_1, y_p)) \rangle$  не является сильно связным, и вершина  $y_1$  принадлежит его начальной компоненте, а вершина  $y_p$  — конечной компоненте, то из леммы 1, предложения 3 и из определений пути  $H$  и  $S_i$  ( $0 \leq i \leq v$ ) следует, что он содержит такой гамильтоновы  $(y_1, y_p)$ -путь  $F$ , что по направлению пути  $F$  вершины  $y_i$  ( $1 \leq i \leq p$ ) с меньшими индексами встречаются раньше, чем вершины  $y_i$  ( $1 \leq i \leq p$ ) с большими индексами. Следовательно,  $C: (P \rightarrow F)$  является гамильтоновым контуром для  $G$ . Пусть для определенности  $C: (z_1 \rightarrow z_n \rightarrow z_{n-1} \rightarrow \dots \rightarrow z_2 \rightarrow z_1)$ , где  $z_1 = x_1$  и  $z_n = y_1$  (заметим, что одни и те же вершины орграфа  $G$  для удобства обозначаем несколькими буквами).

В дальнейшем будем использовать следующие обозначения:

$$\varphi(0) = 0 \quad \text{и} \quad \varphi(k) = \sum_{i=1}^k (m_i + s_i), \quad \text{при} \quad k \in [1, v];$$

$$\psi(k) = \varphi(k-1) + m_k, \quad \text{при} \quad k \in [1, v].$$



Тогда очевидно, что

$$P^i: (z_{\psi(i)} \rightarrow z_{\psi(i)-1} \rightarrow \dots \rightarrow z_{\varphi(i-1)+1}), \quad \text{где } i \in [1, v];$$

$$S_i: (z_{\varphi(i)} \rightarrow z_{\varphi(i)-1} \rightarrow \dots \rightarrow z_{\psi(i)+1}), \quad \text{где } i \in [1, v].$$

Далее для любого  $j \in [0, v]$  пусть (см. рис. 1)

$$A_0 = V(G_1) \quad \text{и} \quad A_j = V(H[y_1, z_{\varphi(j)}]), \quad \text{при } j \geq 1;$$

$$D_v = V(G_s) \quad \text{и} \quad D_j = V(H(z_{\psi(j)+1}, y_p]), \quad \text{при } 1 \leq j \leq v-1;$$

$$D_0 = V(H(y_{s_0}, y_p]).$$

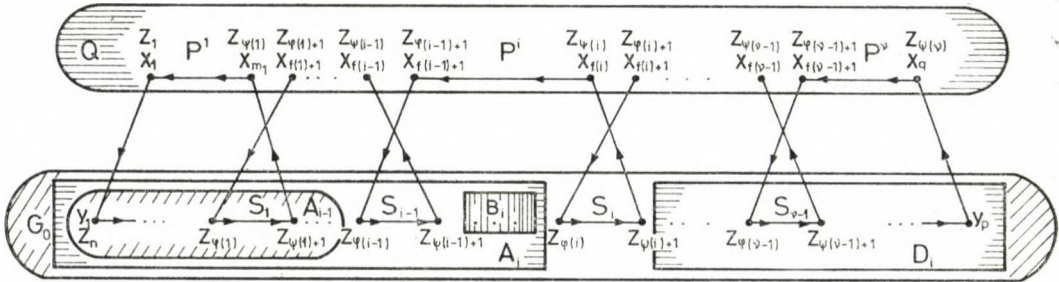


Рис. 1

Из определения подмножеств  $A_j$ ,  $D_j$  и пути  $H$  непосредственно вытекает следующее

**Утверждение 4.** Вершины подмножества  $A_j$  (соответственно  $D_j$ ) составляют подпуть с началом  $y_1$  (с концом  $y_p$ ) для пути  $H$  и

$$V(G_0) = A_j \cup D_j \cup V(S_j); \quad A_j \rightarrow V(G_0) \setminus A_j;$$

$$A_0 \subseteq A_1 \subseteq \dots \subseteq A_v; \quad D_v \subseteq D_{v-1} \subseteq \dots \subseteq D_1 \subseteq D_0.$$

**Утверждение 5.** Для любых  $i \in [1, v-1]$  имеет место

а)  $E(z_{\varphi(i)+1} \rightarrow A_i) = E(D_i \rightarrow z_{\psi(i)}) = \emptyset;$

б)  $L_4 \triangleq E(V(P^{i+1}) \setminus \{z_{\varphi(i)+1}\} \rightarrow A_i \cup V(S_i)) = \emptyset;$

в)  $L_5 \triangleq E(D_i \cup V(S_i) \rightarrow V(P^i) \setminus \{z_{\psi(i)}\}) = \emptyset.$

**Доказательство.** Справедливость утверждения а) непосредственно следует из определения пути  $H$ , подмножеств  $A_i$ ,  $D_i$ , пути  $S_i$ , минимальности числа  $i_i$  и максимальности числа  $i_i$ .

Если  $L_4 \neq \emptyset$  и  $uv \in L_4$ , то из определения пути  $H$  имеем, что путь  $(P[y_p, u] \rightarrow H[v, z_{\varphi(j)+1}] \rightarrow P[z_{\varphi(j)}, z_n])$  принадлежит множеству  $R$  и не проходит через вершину  $z_{\varphi(i)+1}$ , где  $j$  ( $j \geq i$ ) такое минимальное число, что  $v \in A_j \cup V(S_j)$ . Если же  $L_5 \neq \emptyset$  и  $uv \in L_5$ , то путь  $(P[y_p, z_{\varphi(j)+1}] \rightarrow H[z_{\varphi(j)}, u] \rightarrow P[v, z_n])$  принадлежит множеству  $R$  и не проходит через вершину  $z_{\psi(i)}$ , где  $j$  ( $j \geq i$ ) такое максимальное число, что  $u \in D_j \cup V(S_j)$ . Итак, в обоих случаях б) и в) построили

путь, который не проходит через все вершины множества  $V(Q)$ , а это противоречит соотношению (2).

В дальнейшем будем предполагать, что  $G$  не является панциклическим и  $C(G) = \emptyset$ . Отсюда, в частности, имеем, что для некоторого  $m \in [3, n-1]$   $G$  не содержит контура длины  $m$ . Значит для любого  $i \in [1, n]$  имеет место  $z_i z_{i+m-1} \notin E(G)$ , в частности  $z_1 z_m \notin E(G)$ .

**Утверждение 6.** Для некоторого  $i \in [0, v-1]$  число  $m$  удовлетворяет неравенству  $\varphi(i) + 2 \leq m \leq \psi(i+1) + 1$ , т.е.  $z_m \in \{z_{\varphi(i)+2}, z_{\varphi(i)+3}, \dots, z_{\psi(i+1)+1}\}$ .

Действительно, в противном случае имеем  $z_{m-1} \in V(G_0)$ . Следовательно, если  $z_n z_{m-1} \in E(G)$ , то  $C_m: (z_1 \rightarrow z_n \rightarrow C[z_{m-1}, z_1])$ . Если же  $z_n z_{m-1} \notin E(G)$ , то по (1)  $z_{m-1} \in V(G_1)$ . Пусть  $z_t$  та вершина контура  $C$ , который на гамильтоновом пути  $H$  непосредственно предшествует вершине  $z_{\varphi(1)}$ . Тогда, очевидно, что  $A_1 = \{z_n, z_{n-1}, \dots, z_t\}$ ;  $z_{t-1} \notin A_1$  и  $m > t$ . Поэтому по утверждению 4 имеем  $C_m: (C[z_1, z_{n-m+t}] \rightarrow C[z_{t-1}, z_1])$ , а это противоречит предположению  $C_m \not\subset G$ . Утверждение 6 доказано.

Поскольку  $m \geq 3$  и  $z_1 z_m \notin E(G)$ , то из утверждений 1, 5 и 6 следует, что  $E(z_1, z_m) = \emptyset$ .

Пусть для определенности  $\varphi(k) + 2 \leq m \leq \psi(k+1) + 1$ , где  $k \in [0, v-1]$ , и пусть

$$N = A_k \setminus V(C[z_{\varphi(k)}, z_1]) = \{y^1 = y^1, y^2, \dots, y^l\},$$

где элементы множества  $N$  пронумерованы таким образом, что при  $1 \leq i < j \leq l$  по направлению пути  $H$  вершина  $y^j$  встречается после вершины  $y^i$ . Очевидно, что  $(y^1 \rightarrow y^2 \rightarrow \dots \rightarrow y^l) \subset G$ , а из утверждения 4 следует  $N \rightarrow z_{\varphi(k)}$ .

**Утверждение 7.**  $m - \varphi(k) > l$ .

**Доказательство.** Предположим, что  $l \geq m - \varphi(k)$ . Если  $k \geq 1$ , то  $\{z_n, z_{n-1}, \dots, z_{n-(m-\varphi(k))+1}\} \subseteq N$ . Следовательно, поскольку  $N \rightarrow z_{\varphi(k)}$ , то  $C_m: (z_1 \rightarrow z_n \rightarrow z_{n-1} \rightarrow \dots \rightarrow z_{n-(m-\varphi(k))+1} \rightarrow C[z_{\varphi(k)}, z_1])$ . Если же  $k = 0$ , то  $N = V(G_1)$  и  $m - \varphi(0) = m \leq l = s_0$ . Поэтому, согласно лемме 1,  $G_1$  содержит контур длины  $m$ . В обоих случаях получили, что  $C_m \subset G$ , а это противоречит  $C_m \not\subset G$ . Таким образом утверждение 7 доказано.

Из утверждения 7 следует, что

$$(3) \quad \{z_{m-1}, z_{m-2}, \dots, z_{m-l}\} \subseteq V(P^{k+1}).$$

**Утверждение 8.**  $m - s_0 \geq 2$ .

**Доказательство.** Предположим, что  $m - s_0 \leq 1$ . Отсюда, так как  $m \geq 3$ , то  $s_0 \geq 2$ . По лемме 1,  $G_1$  содержит контур любой длины  $i \in [3, s_0]$ . Следовательно,  $m = s_0 + 1$  и, так как  $C_m \not\subset G$ , то из предложения 1 следует

$$d(z_1, V(G_1)) \leq s_0.$$

Далее из  $s_0 \geq 2$  следует, что  $z_{n-s_0+2} \in V(G_1)$  и

$$E(V(G_0) \setminus V(G_1) \rightarrow z_1) = \emptyset,$$

поскольку иначе с помощью (1) имеем  $C_m: (C[z_1, z_{n-s_0+2}] \rightarrow u \rightarrow z_1)$ , где  $u \in$

$\in V(G_0) \setminus V(G_1)$ . Отсюда,  $m_1 \geq 2$  и так как  $E(z_1, z_m) = \emptyset$ , то, пользуясь утверждением 1, получим

$$\begin{aligned} n &\leq d(z_1) = d(z_1, V(Q)) + d(z_1, V(G_1)) + d(z_1, V(G_0) \setminus V(G_1)) \leq \\ &\leq |E(z_1 \rightarrow V(Q) \cup V(G_0) \setminus (V(G_1) \cup \{z_1, z_m\}))| + s_0 + 1 \leq \\ &\leq n - s_0 - 2 + s_0 + 1 \leq n - 1. \end{aligned}$$

Полученное противоречие доказывает справедливость утверждения 8.

Утверждение 9. Для всех  $i \in [1, l]$  имеет место  $E(z_{n-i+1}, z_{m-i}) = \emptyset$ .

Доказательство. Из (3) и утверждения 8 следует, что  $z_{m-i} \in V(P^{k+1}) \setminus \{z_1\}$ ,  $i \in [1, l]$ . Следовательно, поскольку  $z_{n-i+1} \in A_k$ , то согласно утверждениям 3 и 5 имеем  $E(z_{m-i} \rightarrow z_{n-i+1}) = \emptyset$ . Далее, из  $(z_1 \rightarrow z_n \rightarrow z_{n-1} \rightarrow \dots \rightarrow z_{n-i+1}) \subset G$  вытекает, что  $C_m: (z_1 \rightarrow z_n \rightarrow z_{n-1} \rightarrow \dots \rightarrow z_{n-i+1} \rightarrow C[z_{m-i}, z_1])$ , при  $z_{n-i+1} z_{m-i} \in E(G)$ . Следовательно,  $E(z_{n-i+1} \rightarrow z_{m-i}) = \emptyset$ . Итак,  $E(z_{n-i+1}, z_{m-i}) = \emptyset$ . Что и требовалось доказать.

Покажем, что для любого  $j \in [1, l]$  справедливо неравенство

$$(4) \quad d(z_{m-j}, V(S_k)) \leq \begin{cases} s_k + 1, & \text{при } k \geq 1 \text{ и } m-j = \varphi(k) + 1; \\ s_k, & \text{при } k \geq 1 \text{ и } m-j \neq \varphi(k) + 1; \\ s_k - 1, & \text{при } k = 0. \end{cases}$$

Доказательство (4). По (3) имеем  $z_{m-j} \in V(P^{k+1})$ ,  $j \in [1, l]$ . Пусть  $k \geq 1$ . Если  $m-j \neq \varphi(k) + 1$ , то по утверждению 5 имеем  $E(z_{m-j} \rightarrow V(S_k)) = \emptyset$ , и поэтому,  $d(z_{m-j}, V(S_k)) \leq s_k$ . Если же  $m-j = \varphi(k) + 1$ , то из  $N \rightarrow z_{\varphi(k)}$  и  $z_{n-j+1} \in N$  имеем  $z_{n-j+1} z_{\varphi(k)} \in E(G)$ . Следовательно, контур  $(z_n \rightarrow z_{n-1} \rightarrow \dots \rightarrow z_{n-j+1} \rightarrow C[z_{\varphi(k)}, z_n])$  имеет длину  $m-1$ . Поэтому, так как  $S_k$  является подпутем этого контура, то его невозможно расширить с помощью  $z_{\varphi(k)+1}$ . Отсюда, по предложению 2, имеем  $d(z_{m-j}, V(S_k)) \leq s_k + 1$ .

Пусть теперь  $k=0$ . Тогда  $N = V(G_1) = V(S_0)$ , т.е.  $l = s_0$  и, по утверждению 8, имеет место  $m-j \geq 2$ . Следовательно, согласно утверждениям 3 и 9 имеем

$$E(z_{m-j} \rightarrow V(G_1)) = E(z_{n-j+1}, z_{m-j}) = \emptyset.$$

Отсюда  $d(z_{m-j}, V(S_0)) \leq s_0 - 1$ . Итак, неравенство (4) доказано.

Для всех  $i \in [1, \psi(v)]$  обозначим

$$Q_{i,1} = V(Q) \cap V(C(z_i, z_1)); \quad Q_{i,2} = V(Q) \cap V(C(z_{\psi(v)}, z_i)).$$

Нетрудно убедиться, что для любого  $i \in [2, l]$  имеет место

$$(5) \quad E(D_k \rightarrow z_{m-i}) = \emptyset.$$

Действительно, поскольку  $N \rightarrow D_k$ , то в противном случае имеем  $C_m: (z_1 \rightarrow z_n \rightarrow z_{n-1} \rightarrow \dots \rightarrow z_{n-i+2} \rightarrow u \rightarrow C[z_{m-i}, z_1])$ , где  $uz_{m-i} \in E(D_k \rightarrow z_{m-i})$ , а это невозможно.

Утверждение 10. Для всех  $i \in [2, l]$  имеет место  $d(z_{m-i}) = n$  и

$$A_k \cup Q_{m-i,1} \setminus \{z_{n-i+1}\} \rightarrow z_{m-i} \rightarrow D_k \cup Q_{m-i,2}.$$

**Доказательство.** Используя (3) и утверждения 3 и 5 получим  $E(z_{m-i} \rightarrow A_k \setminus \{z_{n-i+1}\}) = \emptyset$ . Следовательно, согласно (4), (5) и утверждениям 1 и 9, имеем

$$n \leq d(z_{m..i}) \leq |E(z_{m-i} \rightarrow D_k \cup Q_{m-i,2})| + |E(Q_{m-i,1} \rightarrow z_{m-i})| + \\ + \begin{cases} |E(A_k \setminus \{z_{n-i+1}\} \rightarrow z_{m-i})| + s_k + 2, & \text{при } k \neq 0; \\ s_k + 1, & \text{при } k = 0, \text{ т.е. } A_0 = V(S_0). \end{cases}$$

Отсюда непосредственно вытекает справедливость утверждения 10.

**Утверждение 11.**  $|E(z_1 \rightarrow V(G_1))| \leq 1$ .

**Доказательство.** Предположим, что  $|E(z_1 \rightarrow V(G_1))| \geq 2$ . Тогда  $s_0 \geq 2$  и пусть  $z_1 y_i \in E(G)$ , где  $2 \leq i \leq s_0$ . Тогда, по утверждению 10,  $V(G_1) \setminus \{y_i\} \rightarrow z_{m-i}$ . Следовательно,  $C_m: (C[z_{m-i}, z_1] \rightarrow C[y_i, y_{(2i-1) \bmod (s_0)}] \rightarrow z_{m-i})$ . Пришли к противоречию, что доказывает утверждение 11.

Теперь покажем, что

$$(6) \quad |E(V(S_1) \rightarrow z_1)| \leq 1.$$

**Доказательство (6).** Предположим, что (6) не верно. Тогда существует такое  $i \in [m_1 + 2, m_1 + s_1]$ , что  $z_i z_1 \in E(G)$ . Отсюда и из утверждения 5в следует  $m_1 = 1$ . Поэтому  $k \geq 1$ . Нетрудно убедиться, что

$$(7) \quad E(D_k \cup \{z_{m_1+1}\} \rightarrow z_{m-i}) = \emptyset.$$

Действительно, в противном случае, поскольку подграф  $\langle \{z_{i-1}, z_{i-2}, \dots, z_2\} \cup D_k \rangle$  содержит такой гамильтонов путь  $(u_1 \rightarrow u_2 \rightarrow \dots \rightarrow u_i)$ , что  $(u_1 \rightarrow u_2 \rightarrow \dots \rightarrow u_{i-2}) = (z_{i-1} \rightarrow z_{i-2} \rightarrow \dots \rightarrow z_2)$  и, так как  $z_n \rightarrow D_k \cup V(S_1)$ , то  $C_m: (z_1 \rightarrow z_n \rightarrow u_{j-i+3} \rightarrow u_{j-i+4} \rightarrow \dots \rightarrow u_j \rightarrow C[z_{m-1}, z_i] \rightarrow z_1)$ , где  $u_j z_{m-1} \in E(G)$  и  $u_j \in D_k \cup \{z_2\}$ , а это невозможно.

Теперь, пользуясь тем, что  $z_2 z_{m-1} \notin E(G)$  аналогично доказательству неравенства (4), получим

$$d(z_{m-1}, V(S_k)) \leq \begin{cases} s_k, & \text{если } m \neq \varphi(k) + 2 \text{ и } k \geq 2 \text{ или } k = 1 \text{ и } m = \varphi(k) + 2; \\ s_k - 1, & \text{если } k = 1 \text{ и } m \neq \varphi(k) + 2; \\ s_k + 1, & \text{если } k \geq 2 \text{ и } m = \varphi(k) + 2. \end{cases}$$

Отсюда, пользуясь утверждениями 2, 3, 5, 9, соотношением (7) и тем, что, если  $z_2 \in A_k$ , то  $E(z_{m-1}, z_2) = \emptyset$ , получим

$$n \leq d(z_{m-1}) = d(z_{m-1}, V(Q)) + |E(A_k \setminus \{z_n, z_2\} \rightarrow z_{m-1})| + \\ + |E(z_{\psi(k)+1} \rightarrow z_{m-1})| + |E(z_{m-1} \rightarrow D_k)| + d(z_{m-1}, V(S_k)) \leq n - 1.$$

Полученное противоречие завершает доказательства неравенства (6).

Для всех  $i \in [1, v]$  через  $B_i$  обозначим  $A_i \cap D_{i-1}$ .

**Утверждение 12.** Если  $B_1 = \emptyset$ , то

а)  $z_1 \rightarrow V(G) \setminus (V(G_1) \cup \{z_1, z_m\})$  и  $d(z_1) = n$ ;

б) Для любого  $i \in [1, k]$ , если  $\psi(i) \leq j_1 < j_2 - 1 \leq \varphi(i)$ , то  $z_{j_2} z_{j_1} \notin E(G)$ .

Доказательство. а) Из утверждений 1, 3 и 5 имеем

$$d(z_1) = |E(z_1 \rightarrow V(Q) \cup D_1 \cup V(S_1) \setminus \{z_1, z_m\})| + \\ + |E(V(S_1) \cup \{z_2\} \rightarrow z_1)| + d(z_1, V(G_1)).$$

Отсюда, поскольку  $d(z_1) \geq n$ , то пользуясь неравенством (6) и утверждениями 1, 5, 11 и  $E(z_1, z_m) = \emptyset$ , получим  $d(z_1) = n$  и  $z_1 \rightarrow V(G) \setminus (V(G_1) \cup \{z_1, z_m\})$ . Что и требовалось доказать

б) Предположим, что для некоторых  $j_1$  и  $j_2$ ,  $\psi(i) \leq j_1 < j_2 - 1 \leq \varphi(i)$ , имеет место  $z_{j_2} z_{j_1} \in E(G)$ . Если  $\varphi(v) - m \geq j_2 - j_1 - 1$ , то с помощью утверждения 12а имеем  $C_m: (z_1 \rightarrow C[z_{m+j_2-j_1-1}, z_{j_2}] \rightarrow C[z_{j_1}, z_1])$ . Пусть  $\varphi(v) - m < j_2 - j_1 - 1$ . Тогда по (1) и утверждению 12а имеет место

$$\{z_1\} \rightarrow \{z_{j_1+1}, z_{j_1+2}, \dots, z_{j_2-1}\} \rightarrow \{z_{\varphi(v)}\}$$

и, поэтому  $C_m: (z_1 \rightarrow C[z_{j_1+\alpha}, z_{j_1+1}] \rightarrow C[z_{\varphi(v)}, z_{j_2}] \rightarrow C[z_{j_1}, z_1])$ , где  $\alpha = j_2 - j_1 - \varphi(v) + m - 1$ , а это невозможно. Итак утверждение 12 доказано.

Утверждение 13. Если  $z_i, z_j \in V(P^1)$  и  $(z_j \rightarrow u_1 \rightarrow u_2 \rightarrow \dots \rightarrow u_i \rightarrow z_i) \subset G$ , где  $i \in [1, v]$ ,  $i < j$  и  $u_1, u_2, \dots, u_i \in B_i$ , то  $j = i + 1$ .

Действительно, в противном случае, из пути  $(C[z_{\psi(v)+1}, z_j] \rightarrow u_1 \rightarrow u_2 \rightarrow \dots \rightarrow u_i \rightarrow C[z_i, z_n])$  получим путь, который принадлежит множеству  $R$  и не проходит через вершину  $z_{i+1}$ , что противоречит соотношению (2).

Нетрудно убедиться, что

$$(8) \quad E(D_k \rightarrow z_{m-2}) = \emptyset,$$

поскольку иначе  $C_m: (z_n \rightarrow u \rightarrow C[z_{m-2}, z_n])$ , где  $uz_{m-2} \in E(D_k \rightarrow z_{m-2})$ .

Утверждение 14. Пусть  $z_{m-2} \in V(P^{k+1})$  и  $m-2 \neq \varphi(k)+1$ . Если для некоторой вершины  $u \in V(G_0)$  имеет место  $E(u, z_{m-2}) = \emptyset$ , то  $d(z_{m-2}) = n$  и

$$Q_{m-2,1} \cup A_k \cup S_k \setminus \{u\} \rightarrow z_{m-2} \rightarrow Q_{m-2,2} \cup D_k \setminus \{u\}.$$

Доказательство. С помощью утверждений 1, 5 и равенства (8) получим

$$d(z_{m-2}) = |E(z_{m-2} \rightarrow D_k \cup Q_{m-2,2} \cup \{z_{m-3}\} \setminus \{u\})| + \\ + |E(Q_{m-2,1} \cup A_k \cup V(S_k) \cup \{z_{m-1}\} \setminus \{u\} \rightarrow z_{m-2})|.$$

Отсюда и из  $d(z_{m-2}) \geq n$  непосредственно следует справедливость утверждения 14.

Заметим, что для любой вершины  $u \notin V(C[z_{m-1}, z_1])$  имеет место

$$(9) \quad |E(z_1 \rightarrow u)| + |E(u \rightarrow z_{m-1})| \leq 1.$$

Действительно, иначе  $C_m: (z_1 \rightarrow u \rightarrow C[z_{m-1}, z_1])$ , что невозможно.

Утверждение 15.  $m > (n+1)/2$ .

Доказательство. Поскольку  $|D_k \cup Q_{m-1,2}| = n - m$ , то для доказательства неравенства  $m > (n+1)/2$  достаточно показать, что

$$(10) \quad z_{m-1} \rightarrow D_k \cup Q_{m-1,2}.$$

Действительно, при справедливости (10)  $G$  содержит контур любой длины  $j \in [2, n-m+1]$ . Значит,  $m > (n+1)/2$ .

Докажем соотношение (10). Если  $l \geq 2$ , то (10) следует из утверждения 10. Поэтому будем предполагать, что  $l=1$ . Тогда  $s_0=l=1$ .

Сначала покажем, что

$$(11) \quad z_1 \rightarrow D_k \setminus \{z_m\}.$$

Допустим, что соотношение (11) неверно, а именно, для некоторой вершины  $v \in D_k \setminus \{z_m\}$  имеет место  $z_1 v \notin E(G)$ . Отсюда и из утверждения 12а вытекает, что  $B_1 \neq \emptyset$ . Следовательно, поскольку  $l=1$ , то  $k=0$ . Значит,  $D_0 = V(G_0) \setminus \{z_n\}$ . Пользуясь утверждениями 1, 5 и  $E(z_1, z_m) = \emptyset$  получим

$$d(z_1) = |E(z_1 \rightarrow V(G) \setminus \{z_1, z_m, v\})| + |E(\{z_2\} \cup A_1 \rightarrow z_1)|.$$

Отсюда и из  $d(z_1) \geq n$  вытекает, что  $|E(A_1 \cup \{z_2\} \rightarrow z_1)| \geq 3$ . Следовательно, для некоторой вершины  $u \in A_1 \setminus \{z_n\} = B_1$  имеет место  $uz_1 \in E(G)$ . Отсюда, поскольку  $z_n u \in E(G)$ , то  $m \geq 4$ .

Докажем, что для любой вершины  $y \in V(G_0)$

$$(12) \quad |E(z_{m-2} \rightarrow y)| + |E(y \rightarrow z_1)| \leq 1.$$

Допустим, что (12) неверно. Тогда для некоторой вершины  $y \in V(G_0)$  имеет место  $z_{m-2}y, yz_1 \in E(G)$ . Отсюда и из утверждений 3 и 5 следует  $y \in B_1$ . Следовательно, по утверждению 13,  $m-2=2$ , т.е.  $m=4$ . Очевидно, что  $E(z_n, z_2) = \emptyset$ , так как в случае  $E(z_n, z_2) \neq \emptyset$  имеем  $z_n z_2 \in E(G)$  и  $C_m: (z_1 \rightarrow z_n \rightarrow z_2 \rightarrow y \rightarrow z_1)$ . Поэтому, по утверждению 14, имеет место  $z_2 \rightarrow V(G) \setminus \{z_n, z_2\}$ . Значит,  $C_4 \subset G$ . Это противоречие доказывает неравенство (12).

Из  $uz_1 \in E(G)$  и (12) имеем  $z_{m-2}u \notin E(G)$ . Значит, по (8)  $E(z_{m-2}, u) = \emptyset$ . Вновь пользуясь утверждением 14 получим  $z_n z_{m-2} \in E(G)$ ,  $d(z_{m-2}) = n$  и

$$(13) \quad z_{m-2} \rightarrow V(G_0) \setminus \{z_n, u\}.$$

Отсюда, согласно (12) имеет место  $E(V(G_0) \setminus \{z_n, u\} \rightarrow z_1) = \emptyset$ . Поэтому, по утверждению 1, имеем

$$d(z_1) = |E(z_1 \rightarrow V(G) \setminus \{z_1, z_m, v\})| + |E(\{z_2, z_n, u\} \rightarrow z_1)|.$$

Следовательно, поскольку  $d(z_1) \geq n$ , то

$$(14) \quad z_1 \rightarrow V(G) \setminus \{z_1, z_m, v\}.$$

Отсюда, в частности, имеем  $z_1 z_{m-1}, z_1 z_{m-2} \in E(G)$ . Далее очевидно, что

$$(15) \quad E(V(G_0) \setminus \{u, z_n\} \rightarrow u) = \emptyset.$$

Действительно, в противном случае из  $uz_1 \in E(G)$  и из (13) с помощью утверждения 13 получим  $m=4$ . Значит  $C_4: (z_1 \rightarrow z_{m-2} \rightarrow w \rightarrow u \rightarrow z_1)$ , где  $w \in V(G_0) \setminus \{u, z_n\}$ , а это невозможно.

В силу соотношения (15) имеем

$$(16) \quad d(u, V(G_0)) = |E(u \rightarrow V(G_0) \setminus \{z_1, u\})| + |E(z_n \rightarrow u)| = p-1.$$



Далее, так как  $uz_1, z_1z_{m-1} \in E(G)$ , то легко заметить, что  $E(V(Q) \setminus \{z_1\} \rightarrow u) = \emptyset$ . Отсюда, пользуясь (9) и учитывая равенство  $E(u, z_{m-2}) = \emptyset$  получим

$$d(u, V(Q)) = |E(u \rightarrow V(Q) \setminus \{z_{m-2}\})| + |E(z_1 \rightarrow u)| \leq q-1.$$

Следовательно, по (16),  $d(u) \leq n-2$ . Поэтому, так как  $d(z_{m-2}) = n$ , то  $d(u) + d(z_{m-2}) \leq 2n-2$ , а это противоречит условию  $(N_0)$ . Полученное противоречие завершает доказательство соотношения (11).

Из (9) и (11) вытекает, что  $E(D_k \setminus \{z_m\} \rightarrow z_{m-1}) = \emptyset$ . Отсюда, пользуясь утверждениями 1, 3, 5, неравенством (4) и учитывая  $E(z_{m-1}, z_n) = \emptyset$  получим: если  $k \geq 1$ , то

$$d(z_{m-1}) = |E(A_k \cup Q_{m-1,1} \setminus \{z_n\} \rightarrow z_{m-1})| + d(z_{m-1}, V(S_k)) + \\ + |E(z_{m-1} \rightarrow Q_{m,2} \cup D_k \cup \{z_{m-2}\} \setminus \{z_m\})| + d(z_{m-1}, z_m);$$

если  $k=0$ , то

$$d(z_{m-1}) = |E(z_{m-1} \rightarrow Q_{m,2} \cup D_k \cup \{z_{m-2}\} \setminus \{z_m\})| + \\ + d(z_{m-1}, z_m) + |E(Q_{m-1,1} \rightarrow z_{m-1})|.$$

Поскольку  $d(z_{m-1}) \leq n$ , то из последних равенств, согласно утверждения 1 и (4) следует, что  $z_{m-1} \rightarrow D_k \cup Q_{m-1,2}$ . Итак, как доказательство соотношения (10) так и доказательство утверждения 15 завершены.

Утверждение 16. Для любого  $i \in [1, k+1]$  имеет место  $B_i = \emptyset$ .

Доказательство. Сначала покажем справедливость утверждения для всех  $i \in [1, k]$ . Поэтому сначала будем предполагать, что  $k \geq 1$ . Заметим, что нужно доказать равенство  $l = s_0$ . Предположим, что  $l \geq s_0 + 1$ . Поскольку  $z_n \rightarrow \{y^2, y^{s_0+1}, y^{s_0+2}, \dots, y^l\}$ , то из (8) и утверждения 3, 5 имеем

$$E(\{y^2, y^{s_0+1}, y^{s_0+2}, \dots, y^l\}, z_{m-2}) = \emptyset.$$

Отсюда и из утверждения 10 вытекает, что  $l=2$ ,  $s_0=1$  (т.е.  $V(G_1) = \{z_n\}$ ), и  $d(z_{m-2}) = n$ . Следовательно, по условию  $(N_0)$ ,  $d(y^2) = n-1$ .

Пусть  $z_1z_{m-1} \in E(G)$ . Тогда путь  $C[z_{\psi(k)}, z_1]$  невозможно расширить с помощью вершины  $y^2$ , поскольку в противном случае имеем  $C_m: (z_1 \rightarrow C[z_{m-1}, z_{\psi(k)+1}] \rightarrow C'[z_{\psi(k)}, z_1])$ , где  $C'[z_{\psi(k)}, z_1]$  расширенный путь. Поэтому, по предположению 2 имеет место

$$(17) \quad d(y^2, V(C[z_{\psi(k)}, z_1])) \leq \begin{cases} \psi(k), & \text{если } z_1y^2 \notin E(G), \\ \psi(k)+1, & \text{если } z_1y^2 \in E(G). \end{cases}$$

Далее, поскольку  $y^2 \in A_k$ , то по утверждению 5 и определению контура  $C$  имеем  $E(y^2 \rightarrow z_n) = \emptyset$ , и

$$E(D_k \cup V(S_k) \cup Q_{\varphi(k),2} \rightarrow y^2) = \emptyset.$$

Отсюда, учитывая неравенство (9) с помощью (17) получим

$$n-1 = d(y^2) = d(y^2, V(C[z_{\psi(k)}, z_1])) + |E(y^2 \rightarrow Q_{\varphi(k), 2} \setminus \{z_{m-1}, z_{m-2}\})| + \\ + d(y^2, V(G_0) \setminus V(C[z_{\psi(k)}, z_1])) + |E(y^2 \rightarrow z_{m-1})| \leq n-2,$$

и приходим к противоречию.

Пусть теперь  $z_1 z_{m-1} \notin E(G)$ . Тогда из утверждения 1 следует, что  $E(z_1, z_{m-1}) = \emptyset$ . Далее пользуясь утверждениями 1, 3 и 5 получим

$$d(z_1) = |E(z_1 \rightarrow V(Q) \cup D_1 \setminus \{z_1, z_2, z_{m-1}, z_m, y^2\})| + \\ + d(z_1, V(S_1)) + d(z_1, \{z_2, z_n, y^2\}).$$

Вновь пользуясь утверждениями 1, 5 и неравенством (6) получим  $z_1 \rightarrow V(G_0) \setminus \{z_m\}$ . Отсюда и из (9) имеем  $E(D_k \setminus \{z_m\} \rightarrow z_{m-1}) = \emptyset$ . Поэтому

$$d(z_{m-1}) = d(z_{m-1}, V(Q) \setminus \{z_1, z_m\}) + |E(A_k \cup V(S_k) \setminus \{y^1, y^2\} \rightarrow z_{m-1})| + \\ + |E(z_{m-1} \rightarrow D_k \setminus \{z_m\})| + d(z_{m-1}, z_m) \leq n-1,$$

а это противоречит условию  $z_{m-1} \in V(Q)$ . Итак доказано, что при всех  $i \in [1, k]$  имеет место  $B_i = \emptyset$ .

Теперь переходим к доказательству  $B_{k+1} = \emptyset$ . Допустим, что  $B_{k+1} \neq \emptyset$ . Выберем такое число  $t$ , что  $|V(C[z_{\varphi(k+1)}, z_t])| = m-1$ . Из  $m > (n+1)/2$  и из предположения  $B_{k+1} \neq \emptyset$  легко получаем, что  $2m-1 > n \geq s_0 + t + m-1$ . Отсюда,  $2 \leq t \leq m-s_0-1$ .

Докажем, что

$$(18) \quad E(z_t, B_{k+1}) = \emptyset.$$

Доказательство (18). Нетрудно убедиться, что

$$L_6 \triangleq E(z_t \rightarrow A_{k+1} \setminus V(C[z_{\varphi(k+1)}, z_t])) = \emptyset,$$

поскольку в противном случае, ввиду  $A_{k+1} \rightarrow z_{\varphi(k+1)}$  имеем  $C_m: (z_t \rightarrow u \rightarrow C[z_{\varphi(k+1)}, z_t])$ , где  $z_t u \in L_6$ , а это невозможно. Отсюда, так как  $B_{k+1} \neq \emptyset$  имеем, что вершина  $z_t \in V(Q)$  и не является концом пути  $P^i$ , где  $i \in [1, k+1]$ . Кроме того из  $L_6 = \emptyset$  следует, что  $E(z_t \rightarrow B_{k+1}) = \emptyset$ .

Остается показать, что  $E(B_{k+1} \rightarrow z_t) = \emptyset$ . Предположим, что  $E(B_{k+1} \rightarrow z_t) \neq \emptyset$ . Тогда для некоторой вершины  $u \in B_{k+1}$  имеет место  $uz_t \in E(G)$ . Отсюда и из утверждения 5в имеем  $z_t \in V(P^{k+1})$ , а из (8), так как  $t \leq m-2$ , следует что  $t \leq m-3$ .

Покажем, что для любой вершины  $w \in B_{k+1}$  имеет место

$$(19) \quad |E(z_{m-2} \rightarrow w)| + |E(w \rightarrow z_t)| \leq 1.$$

Допустим, что (19) неверно, т.е. для некоторой вершины  $w \in B_{k+1}$  имеет место  $z_{m-2}w, wz_t \in E(G)$ . Тогда из  $t \leq m-3$  и из утверждения 13 следует, что  $t = m-3$ . Отсюда  $E(z_n, z_{m-2}) = \emptyset$ , так как в случае  $E(z_n, z_{m-2}) \neq \emptyset$  имеет место  $z_n z_{m-2} \in E(G)$  и, значит,  $C_m: (z_1 \rightarrow z_n \rightarrow z_{m-2} \rightarrow w \rightarrow C[z_{m-3}, z_1])$ . Поэтому, по утверждению 10,  $s_0 = 1$  и, так как  $z_{m-2}$  является внутренней вершиной пути  $P^{k+1}$ , то по утверждению 14, имеет место  $d(z_{m-2}) = n$ . Отсюда пользуясь утверждением 3

и учитывая, что  $E(z_n, \{z_{m-1}, z_{m-2}\}) = \emptyset$ , получим

$$d(z_{m-2}) + d(z_n) = n + |E(z_n \rightarrow V(G) \setminus \{z_n, z_{m-1}, z_{m-2}\})| + |E(z_1 \rightarrow z_n)| \leq 2n - 2,$$

а это противоречит условию  $(N_0)$ , что и доказывает неравенство (19).

Поскольку  $uz_t \in E(G)$ , то из (8) и (19) следует, что  $E(z_{m-2}, u) = \emptyset$ . Следовательно, по утверждению 14 имеет место  $d(z_{m-2}) = n$ ,  $\{z_1, z_n\} \rightarrow z_{m-2}$  и  $z_{m-2} \rightarrow D_k \setminus \{u\}$ . Отсюда, в частности, имеем  $z_{m-2} \rightarrow B_{k+1} \setminus \{u\}$ . Следовательно,

$$(20) \quad |E(B_{k+1} \setminus \{u\} \rightarrow u)| = \emptyset,$$

поскольку в противном случае в силу утверждения 13 имеет место  $m-2 = t+1$  и, значит,  $C_m: (z_1 \rightarrow z_{m-2} \rightarrow v \rightarrow u \rightarrow C[z_t, z_1])$ , где  $vu \in E(B_{k+1} \setminus \{u\} \rightarrow u)$ , а это невозможно.

Так как  $z_n z_{m-2} \in E(G)$ , то пользуясь утверждениям 5 и 13 получим

$$(21) \quad E(Q_{t,2} \rightarrow u) = \emptyset,$$

а из определения подмножеств  $A_i$ ,  $D_i$  и контура  $C$  имеем

$$(22) \quad E(u \rightarrow A_k) = E(V(S_{k+1}) \cup D_{k+1} \rightarrow u) = \emptyset.$$

Далее нетрудно убедиться, что путь  $C[z_t, z_1]$  невозможно расширить с помощью вершины  $u$ , поскольку в противном случае  $C_m: (z_1 \rightarrow z_n \rightarrow z_{m-2} \rightarrow z_{m-3} \rightarrow \dots \rightarrow z_{t+1} \rightarrow C'[z_t, z_1])$ , где  $C'[z_t, z_1]$  расширенный путь. Поэтому по предложению 2 имеет место

$$(23) \quad d(u, V(C[z_t, z_1])) \cong \begin{cases} t+1, & \text{если } z_1 u \in E(G), \\ t, & \text{если } z_1 u \notin E(G). \end{cases}$$

Следовательно, так как для всех  $i \in [1, k]$ ,  $B_i = \emptyset$  и  $d(z_{m-2}) = n$ , то из (9), (20), (21), (22) и (23) получим

$$\begin{aligned} d(z_{m-2}) + d(u) &= n + d(u, V(C[z_t, z_1])) + |E(A_0 \rightarrow u)| + \\ &+ |E(u \rightarrow Q_{t,2} \cup B_{k+1} \cup D_{k+1} \cup V(S_{k+1}) \setminus \{z_{m-2}, z_{m-1}, u\})| + |E(u \rightarrow z_{m-1})| \leq 2n - 2, \end{aligned}$$

а это противоречит условию  $(N_0)$ . Таким образом соотношение (18) доказано.

Покажем, что

$$(24) \quad d(z_t) = n; \quad z_t z_{\varphi(k+1)} \in E(G) \quad \text{и} \quad |B_{k+1}| \leq 1.$$

Пусть вершина  $z_t$  является внутренней вершиной для некоторого пути  $P^i$ , где  $1 \leq i \leq k+1$ . Тогда, поскольку  $B_j = \emptyset$  при всех  $j \in [1, k]$ , то с помощью утверждений 1, 5 и (18) получим

$$\begin{aligned} d(z_t) &= |E(z_t \rightarrow V(S_i) \cup D_i \cup Q_{t,2} \cup \{z_{t-1}\} \setminus B_{k+1})| + \\ &+ |E(A_{i-1} \cup V(S_{i-1}) \cup Q_{t,1} \cup \{z_{t+1}\} \rightarrow z_t)|. \end{aligned}$$

Отсюда, поскольку  $d(z_t) \geq n$ , следует, что  $d(z_t) = n$ ,  $z_t z_{\varphi(k+1)} \in E(G)$  и  $|B_{k+1}| \leq 1$ .

Пусть теперь вершина  $z_t$  не является внутренней вершиной. Тогда  $z_t$  является начальной вершиной некоторого пути  $P^i$ , т.е.  $t = \psi(i)$  ( $1 \leq i \leq k$ ).

Тогда  $k \geq 1$ , и так как  $B_1 = \emptyset$ , то по утверждению 126 имеет место

$$|E(V(S_i) \rightarrow z_{\psi(i)})| \leq 1.$$

Отсюда, пользуясь утверждениями 1, 3, 5 и (18), получим

$$d(z_t) = |E(z_t \rightarrow Q_{t,2} \cup V(S_i) \cup D_i \cup \{z_{t-1}\} \setminus B_{k+1})| + \\ + |E(Q_{t,1} \cup A_{i-1} \cup V(S_{i-1}) \cup \{z_{\psi(i)+1}\} \rightarrow z_t)|.$$

Следовательно, вновь имеем  $d(z_t) = n$ ,  $z_t z_{\varphi(k+1)} \in E(G)$  и  $|B_{k+1}| \leq 1$ . Таким образом (24) доказано.

Пусть для определенности  $B_{k+1} = \{u\}$ . Тогда, поскольку  $E(u, z_t) = \emptyset$ , то из (24) и из условия  $(N_0)$  вытекает, что  $d(u) = n - 1$ . Поскольку контур  $(z_t \rightarrow C[z_{\varphi(k+1)}, z_t])$  имеет длину  $m - 1$ , то его подпуть  $C[z_{\psi(k+1)}, z_{t+1}]$  невозможно расширить с помощью вершины  $u$ . Поэтому по предложению 2 имеет место

$$(25) \quad d(u, V(C[z_{\psi(k+1)}, z_t])) \equiv \begin{cases} \psi(k+1) - t, & \text{если } uz_{\psi(k+1)} \notin E(G), \\ \psi(k+1) - t + 1, & \text{если } uz_{\psi(k+1)} \in E(G). \end{cases}$$

Далее, нетрудно убедиться, что

$$(26) \quad |E(z_{t-1} \rightarrow u)| + |E(u \rightarrow z_{\psi(k+1)})| \leq 1.$$

Действительно, иначе из  $u \rightarrow V(S_{k+1})$  следует, что  $uz_{\varphi(k+1)-1} \in E(G)$  и  $C_m: (z_{t-1} \rightarrow u \rightarrow C[z_{\varphi(k+1)-1}, z_{t-1}])$ , что невозможно.

Докажем, что

$$(27) \quad L_7 \doteq E(\{z_{t+1}, z_{t+2}, \dots, z_{\psi(k+1)}\} \rightarrow u) \neq \emptyset.$$

Доказательство (26). Предположим, что (27) не верно, т.е.  $L_7 = \emptyset$ . Тогда из  $t \leq m - 2$ ,  $E(z_t, u) = \emptyset$  и  $E(u \rightarrow z_{m-2}) = \emptyset$  следует, что  $E(z_{m-2}, u) = \emptyset$ . Следовательно, в силу утверждения 14 имеем  $z_n z_{m-2} \in E(G)$ . Отсюда контур  $(z_1 \rightarrow z_n \rightarrow C[z_{m-2}, z_1])$  имеет длину  $m - 1$  и его подпуть  $C[z_{t-1}, z_1]$  невозможно расширить с помощью вершины  $u$ . Значит, согласно предположению 2, имеет место

$$d(u, V(C[z_{t-1}, z_1])) \equiv \begin{cases} t, & \text{если } z_1 u \in E(G), \\ t - 1, & \text{если } z_1 u \notin E(G). \end{cases}$$

Поскольку для всех  $i \in [1, k]$ ,  $B_i = \emptyset$ , то из (9), утверждения 5, определение контура  $C$  и вышесказанного имеем

$$d(u) = |E(V(G_1) \rightarrow u)| + |E(u \rightarrow D_{k+1} \cup V(S_{k+1}) \cup Q_{t,2} \setminus \{z_{m-2}, z_{m-1}\})| + \\ + d(u, V(C[z_{t-1}, z_1])) + |E(u \rightarrow z_{m-1})| \leq n - 2,$$

а это является противоречием, что и доказывает (27).

С помощью (27) легко убедиться, что

$$(28) \quad E(u \rightarrow \{z_{t-1}\} \cup \{z_{\varphi(k)+1}, z_{\varphi(k)+2}, \dots, z_{t-2}\}) = \emptyset.$$

Действительно, в противном случае из утверждений 5 и 13 следует, что  $z_t \in V(P^{k+1})$  и  $E(\{z_{t+1}, z_{t+2}, \dots, z_{\psi(k+1)}\} \rightarrow u) = \emptyset$ , а это противоречит (27).

Далее, поскольку  $E(z_i, u) = \emptyset$ , то из (25), (26) и (28) имеем

$$(29) \quad d(u, V(C[z_{\psi(k+1)}, z_{t-1}])) \leq \psi(k+1) - t + 1.$$

Теперь с помощью (27) нетрудно убедиться, что для любой вершины  $x \in V(G) \setminus V(C[z_{\psi(k+1)}, z_{t-1}])$  имеет место

$$(30) \quad |E(x, u)| \leq 1.$$

Действительно, если  $z_t \notin V(P^{k+1})$  то (30) непосредственно следует из утверждения 5 и определения контура  $C$ . Если же  $z_t \in V(P^{k+1})$ , то из (27), для некоторого  $i \in [t+1, \psi(k+1)]$  имеет место  $z_i u \in E(G)$ . Отсюда и из  $C(G) = \emptyset$  вытекает, что  $E(u \rightarrow V(S_k) \cup \{z_{\varphi(k)+1}, z_{\varphi(k)+2}, \dots, z_{t-1}\}) = \emptyset$ . Вновь воспользовавшись утверждением 4 и 5, получим (30).

Из (29) и (30) имеем  $d(u) \leq n-2$ , что противоречит  $d(u) = n-1$ . Полученное противоречие завершает доказательство  $B_{k+1} = \emptyset$  и вместе с тем доказательство утверждения 16.

Утверждение 17. Если для некоторых  $1 \leq i < j \leq s_0$  имеет место  $y_i y_j \in E(G_1)$ , то  $j = i+1$ .

Доказательство. Предположим, что  $y_i y_j \in E(G_1)$  и  $1 \leq i < j-1 \leq s_0-1$ . Тогда по утверждению 10 имеет место  $y_{s_0} \rightarrow \{z_{m-1}, z_{m-2}, \dots, z_{m-s_0+1}\}$ . Следовательно,  $C_m: (C[z_1, y_i] \rightarrow C[y_j, y_{s_0}] \rightarrow C[z_{m-s_0+j-i-1}, z_1])$ , что невозможно.

Из утверждений 12 и 16 следует, что

$$(31) \quad z_1 \rightarrow V(G) \setminus (V(G_1) \cup \{z_1, z_m\}).$$

Если рассмотрим орграф  $\bar{G}$ , который получается из  $G$  после переориентации всех дуг, то из определения  $P^i$ ,  $S_i$  и контура  $C$  следует, что вершина  $z_{\psi(v)}$  в орграфе  $\bar{G}$  выполняет ту же роль, что и вершина  $z_1$  в  $G$ . Значит, вершина  $z_{\psi(v)}$  также не смежна с некоторой вершиной  $x_i$  такой, что  $|V(C[z_{\psi(v)}, z_i])| = m$ . Из  $m > (n+1)/2$  следует, что  $i < m$ . Следовательно, применяя утверждения 11, 12б и 16 к вершине  $z_{\psi(v)}$  легко получим, что справедливо следующее

Утверждение 18. Для любого  $i \in [1, v]$  имеет место  $B_i = \emptyset$  и если  $\psi(i) \leq j_1 < j_2 - 1 \leq \varphi(i)$ , то  $z_{j_2} z_{j_1} \notin E(G)$ .

С помощью доказанных утверждений покажем, что если  $C(G) = \emptyset$  и  $G$  не является панциклическим, то  $G \in \Phi_n^m$ . Нужно показать, что  $G$  удовлетворяет условиям I—IV определения множества  $\Phi_n^m$ .

Пусть  $V(G) = \{u_1, u_2, \dots, u_n\}$ , где  $u_i = z_{n+i-s_0}$ , при  $i \in [1, s_0]$  и  $u_i = z_{i-s_0}$ , при  $i \in [s_0+1, n]$ . Тогда очевидно, что  $G$  удовлетворяет условиям I и II. Так как  $G$  не содержит контура длины  $m$ , то справедливость условий III и IV следует из утверждений 1, 3, 5, 11, 17 и 18. Лемма доказана.

Через  $C^m(G)$  обозначим множество контуров длины  $m$ , принадлежащих множеству  $C(G)$ .

Лемма 3. Пусть  $n$ -вершинный ( $n \geq 3$ ) сильно связный орграф  $G$  удовлетворяет условию  $(N_0)$ . Если для некоторого  $n_0 \in [2, n-3]$  имеет место  $C^{n_0}(G) \neq \emptyset$ , то  $G$  является панциклическим или существует  $m \in [n_0+1, n-1]$  такое, что  $C^m(G) \neq \emptyset$ .

Доказательство. Предположим, что утверждение леммы не верно. Тогда  $G$  не является панциклическим и наибольшее  $n_0$ , для которого  $C^{n_0}(G) \neq \emptyset$  меньше, чем  $n-2$ . Отсюда  $n \geq 5$ .

Пусть  $C_{n_0}: (x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_{n_0} \rightarrow x_1) \in C^{n_0}(G)$  и  $G_1, G_2, \dots, G_s$  компоненты подграфа  $Q = \langle V(G) \setminus V(C_{n_0}) \rangle$ . Пусть  $G_0 = \langle V(C_{n_0}) \rangle$  и для любого  $1 \leq i \leq j \leq s$  подграф  $\langle \bigcup_{k=i}^j V(G_k) \rangle$  обозначим через  $G_{i,j}$ .

Из максимальнойности  $n_0$  и предложения I вытекает, что для любой вершины  $y \in V(Q)$  имеет место

$$(32) \quad d(y, V(G_0)) \leq n_0.$$

Следовательно,

$$(33) \quad d(y, V(Q)) \equiv |V(Q)|.$$

Отсюда нетрудно заметить, что подграф  $Q$  является односторонне связным. Следовательно, для любых  $z_1 \in V(G_1)$  и  $z_2 \in V(G_s)$  в подграфе  $Q$  существует  $(z_1, z_2)$ -путь. Так как  $G$  является сильно связным, то  $E(V(G_0) \rightarrow V(G_1)) \neq \emptyset$  и  $E(V(G_s) \rightarrow V(G_0)) \neq \emptyset$ . Отсюда, в частности следует существование  $i, t \in [1, n_0]$  и таких  $v \in V(G_{k_1})$  и  $w \in V(G_{k_2})$ , где  $1 \leq k_1 \leq k_2 \leq s$ , что  $x_i, v, w, x_t \in E(G)$ . Вершины  $x_i$  и  $x_t$  можно выбрать так, чтобы  $x_i \neq x_t$ . Действительно, в противном случае, имеет место  $E(V(Q), V(G_0) \setminus \{x_i\}) = \emptyset$ . Отсюда получим  $d(v) + d(x_{i+1}) \leq 2n-2$ , а это, так как вершины  $v$  и  $x_{i+1}$  несмежны между собой, противоречит условию  $(N_0)$ . Выберем вершины  $x_i, x_t, v$  и  $w$  таким образом, чтобы  $x_i \neq x_t$  и путь  $G_{n_0}[x_i, x_t]$  имел возможно меньшую длину  $\mu+1$ . Тогда  $0 \leq \mu \leq n-2$ , причем никакая вершина  $x_{i+1}, x_{i+2}, \dots, x_{i+\mu}$  не смежна ни с какой вершиной множества  $V(G_{k_1, k_2})$  и

$$(34) \quad E(\{x_{i+1}, x_{i+2}, \dots, x_{i+\mu}\} \rightarrow V(G_{1, k_1})) = E(V(G_{k_2, s}) \rightarrow \{x_{i+1}, x_{i+2}, \dots, x_{i+\mu}\}) = \emptyset.$$

Рассмотрим следующие возможные случаи.

1) В подграфе  $Q$  существует  $(v, w)$ -путь  $P$  такой, что  $V(Q) \setminus V(P) \neq \emptyset$ .

Тогда  $\mu \geq 1$ . Расширим путь  $G_{n_0}[x_i, x_t]$  с помощью вершин множества  $\{x_{i+1}, x_{i+2}, \dots, x_{i+\mu}\}$  насколько это возможно. В результате получим некоторый путь из вершины  $x_i$  в  $x_t$  длины  $n_0 - d - 1$ , где  $1 \leq d \leq \mu$ . Пусть  $\{u_1, u_2, \dots, u_d\} \subseteq \{x_{i+1}, x_{i+2}, \dots, x_{i+\mu}\}$  и вершины  $u_i$  не принадлежат расширенному пути из  $x_i$  в  $x_t$ . Применяя предложения 2 и учитывая (34), получим

$$\begin{aligned} d(u_1) + d(v) &= d(u_1, V(G_0)) + d(v, V(G_0)) + d(u_1, V(Q)) + d(v, V(Q)) \equiv \\ &\equiv n_0 + d - 1 + n_0 - \mu + 1 + |V(G_{1, k_1-1}) \cup V(G_{k_2+1, s})| + \\ &+ |V(G_{1, k_1-1}) \cup V(G_{k_1+1, s})| + 2|V(G_{k_1})| - 2 \leq 2n-2, \end{aligned}$$

а это противоречит условию  $(N_0)$ .

2) В подграфе  $Q$  любой  $(v, w)$ -путь проходит через все вершины множества  $V(Q)$ .



Пусть  $P: (v=y_1 \rightarrow y_2 \rightarrow \dots \rightarrow y_{n-n_0}=w) \subseteq Q$ . Очевидно, что при всех  $1 \leq i < j-1 \leq n-n_0-1$  имеет место  $y_i y_j \notin E(G)$ . Отсюда следует, что

$$d(y_i, V(Q)) \leq \begin{cases} n-n_0+1, & \text{если } i \notin \{1, n-n_0\}; \\ n-n_0, & \text{если } i \in \{1, n-n_0\}. \end{cases}$$

Значит, поскольку  $d(y_i, V(G_0)) \leq n_0$  и  $d(y_i) \leq n$ , то

$$(35) \quad \{y_2, y_3, \dots, y_{n-n_0}\} \rightarrow y_1; \quad y_{n-n_0} \rightarrow \{y_1, y_2, \dots, y_{n-n_0-1}\}$$

и  $d(y_1, V(G_0)) = d(y_{n-n_0}, V(G_0)) = n_0$ . Следовательно, из непанцикличности  $G$  следует  $E(y_1 \rightarrow V(G_0)) \neq \emptyset$  и существует такие вершины  $x_{i_1}, x_{i_1} \in V(G_0)$ , что  $x_{i_1} y_1, y_1 x_{i_1} \in E(G_0)$ . Так как  $n \geq 5$ , то вершины  $x_{i_1}$  и  $x_{i_1}$  можно выбрать таким образом, чтобы  $x_{i_1} \neq x_{i_1}$ , и подпуть  $G_{n_0}[x_{i_1}, x_{i_1}]$  имел возможно меньшую длину  $\mu_1 + 1$ . Из минимальности  $\mu_1$  имеем  $E(y_1, \{x_{i_1+1}, x_{i_1+2}, \dots, x_{i_1+\mu_1}\}) = \emptyset$ . Отсюда с помощью предложения 2 получим  $\mu_1 = 1$ . С другой стороны, из  $n-n_0 \geq 3$  и (35) следует, что вершина  $x_{i_1+1}$  не смежна ни с какой вершиной  $y_1, y_2, \dots, y_{n-n_0-1}$  и  $x_{i_1+1} y_{n-n_0} \notin E(G)$ . Значит,  $d(x_{i_1+1}, V(Q)) \leq 1$ . Поэтому  $d(y_1) + d(x_{i_1+1}) \leq 2n-2$ , что противоречит условию  $(N_0)$ . Полученное противоречие завершает доказательство леммы 3.

Пусть  $G$  орграф и  $x \in V(G)$ . Обозначим

$$O(x) = \{y \in V(G) / xy \in E(G)\}; \quad I(x) = \{y \in V(G) / yx \in E(G)\}.$$

Предложение 4 ([5]). Пусть сильно связный  $n$ -вершинный орграф  $G$  содержит контур  $C_{n-2}: (x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_{n-2} \rightarrow x_1)$  и пусть  $y_1, y_2 \in V(G) \setminus V(C_{n-2})$ . Если  $d(y_1) \geq n$ ,  $d(y_2) \geq n$  и  $C_{n-1} \not\subseteq G$ , то  $n$  четное и нумерации вершин контура  $C_{n-2}$  можно выбрать так, чтобы

$$O(y_1) = I(y_1) = \{x_1, x_3, x_5, \dots, x_{n-3}, y_2\};$$

$$O(y_2) = I(y_2) = \{x_2, x_4, x_6, \dots, x_{n-2}, y_1\}.$$

Лемма 4. Пусть  $n$ -вершинный ( $n \geq 4$ ) сильно связный орграф  $G$  удовлетворяет условию  $(N_0)$ . Если  $C_{n-2}(G) \neq \emptyset$  и  $G$  не содержит контура длины  $n-1$ , то  $n$  четное и  $G \in \{\bar{K}_{n/2, n/2}, \bar{K}_{n/2, n/2} \setminus \{e\}\}$ , где  $e \in E(\bar{K}_{n/2, n/2})$ .

Доказательство. Пусть  $C_{n-2}: (x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_{n-2} \rightarrow x_1) \in C_{n-2}(G)$  и  $y_1, y_2 \in V(G) \setminus V(C_{n-2})$ . Так как  $d(y_1) \geq n$ , и  $d(y_2) \geq n$ , то по предложению 4  $n$  четное и

$$(36) \quad O(y_1) = I(y_1) = \{x_1, x_3, \dots, x_{n-3}, y_2\},$$

$$O(y_2) = I(y_2) = \{x_2, x_4, \dots, x_{n-2}, y_1\}.$$

Пусть  $x_i, x_j \in \{x_2, x_4, \dots, x_{n-2}\}$  и  $x_i x_j \in E(G)$ . Если  $|\{x_{i+1}, x_{i+2}, \dots, x_{j-1}\}| \geq 3$ , то  $C_{n-1}: (x_{j+1} \rightarrow x_{j+2} \rightarrow \dots \rightarrow x_i \rightarrow x_j \rightarrow y_2 \rightarrow x_{i+2} \rightarrow x_{i+3} \rightarrow \dots \rightarrow x_{j-1} \rightarrow y_1 \rightarrow x_{j+1})$ , а если же  $|\{x_{i+1}, x_{i+2}, \dots, x_{j-1}\}| = 1$ , то  $C_{n-1}: (x_i \rightarrow x_j \rightarrow y_2 \rightarrow y_1 \rightarrow x_{j+1} \rightarrow x_{j+2} \rightarrow \dots \rightarrow x_i)$ . Таким образом в обоих случаях  $C_{n-1} \subset G$ , а это невозможно. Значит  $E(\langle\langle\{x_2, x_4, \dots, x_{n-2}\}\rangle\rangle) = \emptyset$ . Аналогичным образом, доказывается  $E(\langle\langle\{x_1, x_3, \dots, x_{n-3}\}\rangle\rangle) = \emptyset$ .

Отсюда и из (36) имеем

$$E(\langle \{x_1, x_3, \dots, x_{n-3}, y_2\} \rangle) = E(\langle \{x_2, x_4, \dots, x_{n-2}, y_1\} \rangle) = \emptyset.$$

Поэтому из условия  $(N_0)$  непосредственно вытекает, что  $G \in \{\bar{K}_{n/2, n/2}, \bar{K}_{n/2, n/2} \setminus \{e\}\}$ .

Лемма 4 доказана.

**Теорема.** Пусть  $n$ -вершинный  $(n \geq 2)$  сильно связный произвольный орграф  $G$  удовлетворяет условию  $(N_0)$ . Тогда справедливо хотя бы одно из следующих утверждений.

а)  $G$  является панциклическим;

б)  $n$  четное и  $G \in \{\bar{K}_{n/2, n/2}, \bar{K}_{n/2, n/2} \setminus \{e\}\}$ , где  $e \in E(\bar{K}_{n/2, n/2})$ .

в) для некоторого  $n-1 \equiv m \pmod{(n+1)/2}$  имеет место  $G \in \Phi_n^m$ .

**Доказательство.** Нетрудно проверить справедливость теоремы для орграфов с не более 4 вершинами. Предположим, что теорема верна для орграфов с менее  $n$  ( $n \geq 5$ ) вершинами и пусть  $n$ -вершинный орграф  $G$  удовлетворяет условиям теоремы.

Если  $C^{n-1}(G) \neq \emptyset$ , то из предложения 1 следует, что  $G$  является панциклическим. Поэтому будем предполагать, что  $C^{n-1}(G) = \emptyset$ . Из леммы 2,3 следует, что  $G$  является панциклическим или  $G \in \Phi_n^m$  или  $C^{n-2}(G) \neq \emptyset$ . Значит, нужно рассмотреть только случай  $C^{n-2}(G) \neq \emptyset$ . Если  $C_{n-1} \not\subset G$ , то справедливость теоремы следует из леммы 4. Поэтому предположим, что  $C_{n-1} \subset G$ .

Пусть  $C_{n-2}: (x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_{n-2} \rightarrow x_1) \in C^{n-2}(G)$ ;  $G_0 = \langle V(C_{n-2}) \rangle$  и  $y_1, y_2 \notin V(G_0)$ . Поскольку  $C^{n-1}(G) = \emptyset$ , то из предложения 1 следует, что  $y_1 y_2, y_2 y_1 \in E(G)$  и для любого  $i \in [1, 2]$  и  $j \in [1, n-2]$  имеет место  $d(y_i) = n$ ;  $d(y, V(G_0)) = n-2$  и

$$(37) \quad |E(x_j \rightarrow y_i)| + |E(y_i \rightarrow x_{j+1})| \leq 1.$$

Если для некоторого  $i \in [1, 2]$  имеет место  $E(y_i \rightarrow V(G_0)) = \emptyset$  или  $E(V(G_0) \rightarrow y_i) = \emptyset$ , то из  $d(y_i, V(G_0)) = n-2$  и сильной связности  $G$  легко получим, что  $G$  — панциклический. Поэтому будем предполагать, что для любого  $i \in [1, 2]$  имеет место

$$(38) \quad E(y_i \rightarrow V(G_0)) \neq \emptyset; \quad E(V(G_0) \rightarrow y_i) \neq \emptyset.$$

Если для некоторых  $i \in [1, n-2]$  и  $j \in [1, 2]$  имеет место  $x_i y_j, y_j x_{i+2} \in E(G)$ , то  $|E(x_{i+1}, y_{i+1})| = 2$  (индексы вершины  $y$  берутся по  $\text{mod}(2)$ ), так как в противном случае для контура  $C'_{n-2}: (x_i \rightarrow y_j \rightarrow x_{i+2} \rightarrow x_{i+3} \rightarrow \dots \rightarrow x_{i-1} \rightarrow x_i)$  имеет место  $d(y_{j+1}, V(C'_{n-2})) \geq n-1$  и, значит, в силу предложения 1  $G$  является панциклическим.

Покажем, что для любой вершины  $x_i$  ( $1 \leq i \leq n-2$ ) имеет место

$$(39) \quad |E(x_i, \{y_1, y_2\})| \leq 2.$$

Предположим, что для некоторого  $i \in [1, n-2]$  имеет место

$$|E(x_i, \{y_1, y_2\})| \geq 3.$$

Тогда не нарушая общности можно предполагать, что  $|E(x_i, y_1)| = 2$  и  $x_i y_2 \in E(G)$ . С помощью (38) нетрудно убедиться в существовании такого

$i \in [2, n-3]$ , что  $y_1 x_{i+t} \in E(G)$  и

$$(40) \quad E(y_1, x_{i+t-1}) = E(y_1 \rightarrow \{x_{i+1}, x_{i+2}, \dots, x_{i+t-1}\}) = \emptyset.$$

Действительно, в противном случае по (37) имеем  $E(y_1 \rightarrow V(G_0) \setminus \{x_i\}) = \emptyset$ . Отсюда, так как  $d(y_1, V(G_0)) = n-2$ , то  $E(y_1, x_{i-1}) = \emptyset$  и  $\{x_{i+1}, x_{i+2}, \dots, x_{i-2}\} \rightarrow y_1$ . Следовательно,  $C_m: (y_1 \rightarrow x_i \rightarrow x_{i+1} \rightarrow \dots \rightarrow x_{i+m-2} \rightarrow y_1)$  для всех  $m \in [3, n-2]$ . Поэтому, так как  $C_{n-1}, C_n \subset G$ , то  $G$  — панциклический.

Из  $C^{n-1}(G) = \emptyset$  с помощью (40) и предложения 2 легко убедиться, что  $x_{i+t-2} y_1 \in E(G)$ . Отсюда и из  $y_1 x_{i+t} \in E(G)$  имеем, что  $|E(y_2, x_{i+t-1})| = 2$ . Следовательно, поскольку  $x_i y_2 \in E(G)$ , то для некоторой вершины  $x_j \in \{x_{i+1}, x_{i+2}, \dots, x_{i+t-1}\}$  имеет место  $x_{j-1} y_2, y_2 x_{j+1} \in E(G)$ . Поэтому  $|E(x_j, y_1)| = 2$ , а это противоречит соотношению (40), что и завершает доказательство неравенства (39).

Из (39) следует, что орграф  $G_0$  удовлетворяет условиям теоремы. Согласно индуктивному предположению  $G_0$  удовлетворяет утверждению теоремы. Так как  $C_{n-1}, C_n \subset G$ , то из панцикличности  $G_0$  следует панцикличность  $G$ . Поэтому будем предполагать, что  $G_0 \in \Phi_{n-2}^{m_1}$ , где  $n-3 \cong m_1 > (n-1)/2$  или  $(n-2)$ -четное и  $G_0 \in \{\bar{K}_{(n-2)/2, (n-2)/2}, \bar{K}_{(n-2)/2, (n-2)/2} \setminus \{e\}\}$ . Покажем, что в обоих случаях  $G$  является панциклическим.

*Случай 1.*  $G_0 \in \{\bar{K}_{(n-2)/2, (n-2)/2}, \bar{K}_{(n-2)/2, (n-2)/2} \setminus \{e\}\}$ .

Пусть  $A \cup B = V(G_0)$ ,  $A \cap B = \emptyset$  и  $E(\langle A \rangle) = E(\langle B \rangle) = \emptyset$ .

Поскольку  $C_{n-1} \subset G$  и имеет место (38), то не нарушая общности можно предполагать, что  $u y_1 \in E(A \rightarrow y_1)$  и  $y_1 v \in E(y_1 \rightarrow B)$ .

Если  $uv \in E(G_0)$ , то дуга  $uv$  в  $G_0$  находится на контурах длины 2, 4, ...,  $n-4$ . С помощью этих контуров и вершины  $y_1$  получим контуры всех длин  $i \in [2, n-3]$ . Значит, так как  $C_{n-2}, C_{n-1}, C_n \subset G$ , то  $G$  является панциклическим.

Пусть теперь  $uv \notin E(G_0)$ . Тогда для любых вершин  $u_1 \in A$  и  $v_1 \in B$ , где  $\{u, v\} \neq \{u_1, v_1\}$  имеет место  $|E(u_1, v_1)| = 2$ . Учитывая предыдущий случай можем предполагать, что

$$E(y_1 \rightarrow B \setminus \{v\}) = E(A \setminus \{u\} \rightarrow y_1) = \emptyset.$$

Следовательно,  $E(y_1, B \setminus \{v\}) = \emptyset$  или  $E(y_1, A \setminus \{u\}) = \emptyset$ , так как иначе для некоторых вершин имеет место рассмотренный нами случай. Отсюда следует, что  $n=6$ . Поэтому, так как  $C_4, C_5, C_6 \subset G$  и  $C_3: (y_1 \rightarrow v \rightarrow u \rightarrow y_1)$ , то  $G$  — панциклический.

*Случай 2.*  $G_0 \in \Phi_{n-2}^{m_1}$ , где  $n-3 > m_1 > (n-1)/2$ .

Чтобы доказать панцикличность  $G$ , достаточно показать, что  $C_{m_1} \subset G$ . Пусть  $C'_{n-2}: (x_{n-2} \rightarrow x_{n-3} \rightarrow \dots \rightarrow x_2 \rightarrow x_1 \rightarrow x_{n-2}) \subset G_0$ . Тогда для всех  $j \in [2, n-2] \setminus \{m_1\}$  имеет место  $x_1 x_j \in E(G)$  и  $E(x_1, x_{m_1}) = \emptyset$ . Кроме того, вершина  $x_{m_1}$  смежна из всех вершин  $x_2, x_3, \dots, x_{m_1-1}$ . Из (37) и (38) следует, что для некоторого  $i \in [1, n-2]$  имеет место  $y_1 x_i, x_{i+2} y_1 \in E(G)$  и  $E(y_1, x_{i+1}) = \emptyset$ .

2.1)  $i+1 \notin \{1, m_1, n-2\}$ .

Тогда  $x_1 x_{i+1} \in E(G)$ . Очевидно, что  $x_{i+1} x_{n-2} \notin E(G)$ , так как в противном случае  $C_{n-1}: (x_1 \rightarrow x_{i+1} \rightarrow x_{n-2} \rightarrow x_{n-3} \rightarrow \dots \rightarrow x_{i+2} \rightarrow y_1 \rightarrow x_i \rightarrow \dots \rightarrow x_1) \in C^{n-1}(G)$ , что не-

возможно. Поэтому, с помощью условия  $(N_0)$  нетрудно заметить, что  $i+1 \neq n-3$ , так как в случае  $i-1=n-3$  из  $x_{i+1}x_{n-2} \notin E(G)$  следует, что  $m_1 \leq n-4$ ;  $E(x_{n-2}, x_{n-m_1-1}) = \emptyset$  и

$$d(x_{n-2}, V(G_0)) + d(x_{n-m_1-1}, V(G_0)) \leq 2n-6,$$

а это противоречит условию  $(N_0)$ . Значит,  $E(x_{n-2}, x_{i+1}) = \emptyset$ . Отсюда имеем  $i+1 = n-m_1-1$ . Поскольку  $m_1 > n-m_1-1$  и  $x_{j_1}x_{n-m_1-1}, x_{n-m_1-1}x_{j_2} \in E(G)$  при всех  $1 \leq j_1 \leq n-m_1-2$  и при всех  $n-m_1 \leq j_2 \leq n-3$  соответственно, то  $C_{m_1}: (x_2 \rightarrow x_{i+1} \rightarrow x_{m_1} \rightarrow x_{m_1+1} \rightarrow \dots \rightarrow x_{i+2} \rightarrow y_1 \rightarrow x_i \rightarrow \dots \rightarrow x_2)$  или  $C_{m_1}: (x_1 \rightarrow x_2 \rightarrow x_{m_1-1} \rightarrow \dots \rightarrow x_3 \rightarrow y_1 \rightarrow x_1)$ , соответственно для  $3 \leq i+1 \leq m_1-1$  и  $i+1=2$ .

2.2)  $i+1 \in \{1, m_1, n-2\}$ .

Если  $i+1 \in \{1, m_1\}$ , то  $C_{m_1}: (x_3 \rightarrow x_{m_1} \rightarrow x_{m_1+1} \rightarrow y_1 \rightarrow x_{m_1-1} \rightarrow x_{m_1-2} \rightarrow \dots \rightarrow x_3)$  или  $C_{m_1}: (x_1 \rightarrow x_{m_1-2} \rightarrow \dots \rightarrow x_2 \rightarrow y_1 \rightarrow y_2 \rightarrow x_1)$ , соответственно для  $i+1=m_1$  и  $i+1=1$ .

Пусть,  $i+1=n-2$ . Тогда очевидно, что  $x_{n-2}y_2, y_2y_1, x_{n-m_1+1}x_{n-2} \in E(G)$  и  $n-m_1+1 \leq n-3$ . Поэтому,  $C_{m_1}: (x_{n-2} \rightarrow y_2 \rightarrow y_1 \rightarrow x_{n-3} \rightarrow \dots \rightarrow x_{n-m_1+1} \rightarrow x_{n-2})$ .

Таким образом, всевозможные случаи рассмотрены. Теорема доказана.

Из теоремы следуют некоторые из ранее известных результатов, относящихся к гамильтоновости и панциклическости графов и орграфов ([8], [9], [10], [11], [4], [5], [6], [12], [13] и [2]). Приведем некоторые из них.

Следствие 1 ([12]). Если  $n$ -вершинный ( $n \geq 3$ ) сильно связный орграф  $G$  удовлетворяет условию  $(N_2)$ , то  $G$ —панциклический.

Следствие 2 ([6]). Если  $n$ -вершинный ( $n \geq 3$ ) сильно связный орграф  $G$  удовлетворяет условию  $(N_1)$ , то  $G$ —панциклический или  $n$  четное и  $G = \vec{K}_{n/2, n/2}$ .

Следствие 3 ([13]). Если в  $n$ -вершинном графе  $G$  сумма степеней любых двух несмежных различных вершин не меньше  $n$ , то  $G$ —панциклический или  $n$  четное и  $G = K_{n/2, n/2}$ .

Следствие 4 ([4]). Если  $n$ -вершинный ( $n \geq 2$ ) сильно связный орграф  $G$  удовлетворяет условию  $(N_0)$ , то  $G$ —гамильтонов.

Следствие 5 ([2]). Пусть  $n$ -вершинный ( $n \geq 3$ ) сильно связный орграф  $G$  удовлетворяет условию  $(N_0)$ . Тогда  $G$ —гамильтонов и

1) если для некоторого  $m \in [2, n-1]$  имеет место  $C_m \not\subset G$ , то  $G$  содержит контур любой длины  $k \in [2, n] \setminus \{m\}$  или

2)  $n$  четное и  $G \in \{\vec{K}_{n/2, n/2}, \vec{K}_{n/2, n/2} \setminus \{e\}\}$ .

## ЛИТЕРАТУРА

- [1] Харари, Ф., *Теория графов*, Мир, Москва, 1973. MR 49 #10586.
- [2] Дарбинян, С. Х., О панциклических орграфах, ВЦ АН Арм. ССР и ЕрГУ, Ереван, 1979, препринт.
- [3] CHVÁTAL, V., New directions in Hamiltonian Graph Theory, in *New directions in Graph Theory*, Proceedings of the 1971 Ann Arbor Graph Theory Conference, ed. by F. Harary, New York, Academic Press, 1973, 65—95. MR 49 #4821.
- [4] MEYNIEL, M., Une condition suffisante d'existence d'un circuit hamiltonien dans un graphe orienté, *J. Combinatorial Theory Ser. B* 14 (1973), 137—147. MR 47 #6546.

- [5] HÄGGKVIST, R. and THOMASSEN, C., On pancyclic digraphs, *J. Combinatorial Theory Ser. B* **20** (1976) 20—40. *MR* **52** # 10481.
- [6] THOMASSEN, C., An Ore-type condition implying a digraph to be pancyclic, *Discrete Math.* **19** (1977), 85—92. *MR* **58** # 21776.
- [7] GOLDBERG, M. and MOON, J. W., Cycles in  $k$ -strong tournaments, *Pacific J. Math.* **40** (1972), 89—96. *MR* **46** # 3363.
- [8] DIRAC, G. A., Some theorems on abstract graphs, *Proc. London Math. Soc.* **2** (1952) 69—81. *MR* **13**-856.
- [9] ORE, O., Note on Hamilton circuits, *Amer. Math. Monthly* **67** (1960), 55. *MR* **22** # 9454.
- [10] GHOUILA-HOURI, A., Une condition suffisante d'existence d'un circuit hamiltonien, *C. R. Acad. Sci. Paris* **251** (1960), 495—497. *MR* **22** # 5590.
- [11] WOODALL, D. R., Sufficient conditions for circuits in graphs, *Proc. London Math. Soc.* **24** (1972), 739—755. *MR* **47** # 6549.
- [12] OVERBECK-LARISCH, M., A theorem on pancyclic oriented graphs, *J. Combinatorial Theory Ser. B* **23** (1977), 168—173. *MR* **57** # 2976.
- [13] BONDY, J. A., Pancyclic graphs I, *J. Combinatorial Theory Ser. B* **11** (1971), 80—84. *MR* **44** # 2642.

(Поступила 22-ого декабря 1980 г.)

ВЫЧИСЛИТЕЛЬНЫЙ ЦЕНТР  
АКАДЕМИИ НАУК АРМЯНСКОЙ ССР И  
ГОСУДАРСТВЕННОГО УНИВЕРСИТЕТА  
УЛ. П. СЕВАКА 1  
SU—375044 ЕРЕВАН 44  
SOVIET UNION





# THE EXPECTED HEIGHT OF PATHS FOR SEVERAL NOTIONS OF HEIGHT

WOLFGANG PANNY and HELMUT PRODINGER

## Abstract

In this paper lattice paths with two directions are considered. Several notions of height are introduced, namely the maximal deviation, the maximal span and the onesided height. Assuming all paths of length  $n$  to be equally likely, exact enumeration formulae for the expected height and their asymptotic equivalents are derived.

## 1. Introduction

This paper deals with the *expected height of lattice paths* for several notions of height.

A *path of length  $n$*  is a sequence of integers  $a_0, a_1, \dots, a_n$  with  $|a_{i+1} - a_i| = 1$ ,  $0 \leq i < n$ .

Let us just review some previously known results:

If all paths  $a_0, a_1, \dots, a_{2n}$  with  $a_0 = a_{2n} = 0$ ,  $a_i \geq 0$  are assumed to be equally likely and the height of the path is defined to be  $\max \{a_i | 0 \leq i \leq 2n\}$ , then the expected height is

$$(1.1) \quad \sqrt{\pi n} - \frac{3}{2} + O(n^{-1/2+\varepsilon}) \quad \text{for } \varepsilon > 0 \text{ and } n \rightarrow \infty.$$

This result is due to De Bruijn, Knuth and Rice [2] and was stated not in terms of paths but in terms of *planted plane trees*. Such a path can also be considered as a *Dyckword* of length  $2n$ , if the  $i$ -th letter of the word is an opening (closing) bracket for  $a_i - a_{i-1} = 1$  ( $-1$ ).

A word  $u$  is said to be a *prefix* of a word  $w$  iff there exists a word  $v$  with  $uv = w$ . So a prefix of a Dyckword can be considered as a path  $a_0, a_1, \dots, a_n$  with  $a_0 = 0$  and  $a_i \geq 0$ . Assuming all such paths to be equally likely and defining the height again by  $\max \{a_i | 0 \leq i \leq n\}$ , the expected height is

$$(1.2) \quad (\log 2) \sqrt{2\pi n} - \frac{3}{2} + O(n^{-1/2}) \quad \text{for } n \rightarrow \infty.$$

This result is due to Kemp [6].

Regarding paths where a *third direction* is allowed ( $a_{i+1} - a_i = 0$ ), see [10] and [9].

For a fairly exhaustive list of references of related problems see [7].

In Section 2 we consider paths  $a_0, \dots, a_n$  with  $a_0=0$ . The notion of height is defined by  $\max \{|a_i| \mid 0 \leq i \leq n\}$  and is called *maximal deviation*. This class of paths corresponds to the prefixes of the Dyck language, except for the condition  $a_i \geq 0$ , which need not be fulfilled now. We prove that the expected maximal deviation is given by

$$(1.3) \quad \sqrt{\frac{n\pi}{2}} - \frac{1}{2} + O(n^{-1/2+\varepsilon}) \quad \text{for } \varepsilon > 0 \quad \text{and } n \rightarrow \infty.$$

Though the result resembles the previous ones, it cannot be concluded from them as a corollary; to handle the problem we had to use a new technique: Unlike in the former problems, we had not to approximate one single binomial coefficient by the exponential function but a sum of binomial coefficients by the error function. Thus the *Mellin transform* of the error function and the inversion formula of the Mellin transform come into play.

In Section 3 the same family of paths is considered, but the height is now  $\max \{a_i - a_j \mid 0 \leq i, j \leq n\}$ , called *maximal span*. We prove that the expected maximal span is given by

$$(1.4) \quad \sqrt{\frac{8n}{\pi}} - 1 + O(n^{-1/2}) \quad \text{for } n \rightarrow \infty.$$

Though the explicit enumeration formulae are even more complicated than in Section 2, the derivation of the asymptotic formula (1.4) is (due to lucky circumstances) quite elementary.

We would like to mention that the maximal span preserves a way to give a meaningful notion of height not only for Dyckwords of prefixes or Dyckwords but *for all words over a two-letter-alphabet*!

In Section 4 we consider paths  $a_0, a_1, \dots, a_{2n}$  with  $a_0=a_{2n}=0$  and the *onesided height* defined by  $\max \{a_i \mid 0 \leq i \leq 2n\}$ . We prove that the expected onesided height is given by

$$(1.5) \quad \frac{1}{2} \sqrt{n\pi} - \frac{1}{2} + O(n^{-1/2}) \quad \text{for } n \rightarrow \infty.$$

Again the derivation of the asymptotic formula is — in a certain sense — elementary. The following interpretation can be given: Suppose there are two players  $A$  and  $B$  each one having  $n$  cards out of the set  $\{1, \dots, 2n\}$ .  $A$  always leads a card and  $B$  follows. A player makes a trick whenever his number is the greater one. The interesting parameter is *the number of tricks that player  $A$  can make*, if  $B$  plays his optimal strategy. The partitioning of the cards can be seen as a path by defining  $a_0=0$  and  $a_i - a_{i-1} = 1$  ( $-1$ ) iff player  $B$  ( $A$ ) has the card  $i$ . For example,  $A$  has cards 2, 4, 5, 7 means the path 0, 1, 0, 1, 0,  $-1$ , 0,  $-1$ , 0. It is not hard to see that the number of tricks that player  $A$  can make is just the onesided height of the corresponding path!

Section 5 is devoted to some concluding remarks.

In the following sections the *terminus* path is used according to the introduction.

We would like to point out that one notation may have different meanings in different sections.

## 2. The expected maximal deviation

Let  $\psi_{h,l}(z)$  be the generating function where the coefficient of  $z^n$  gives the number of paths from  $(0,0)$  to  $(n,l)$  with maximal deviation  $\leq h$ .

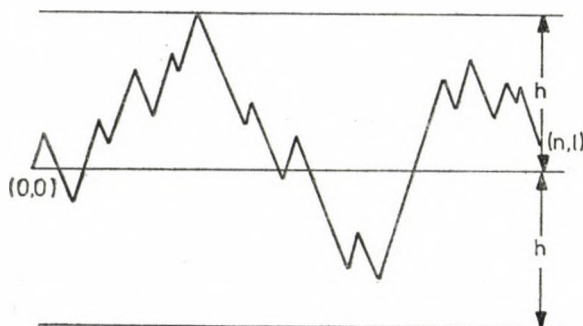


Fig. 1

THEOREM 2.1. With  $z=v/(1+v^2)$  we have

$$(2.1) \quad \psi_{h,l}(z) = \frac{1+v^2}{1-v^2} v^l \frac{1-v^{2(h+1-l)}}{1+v^{2h+2}}.$$

PROOF. Regarding the last step of the path we find the following system of linear recurrences for the generating functions  $\psi_{h,l}$  ( $|l| \leq h$ ) which can be expressed in matrix form:

$$(2.2) \quad \begin{bmatrix} 1 & -z & & & & \\ -z & 1 & -z & & & \\ & & \ddots & \ddots & \ddots & \\ & & & -z & 1 & -z \\ & & & & \ddots & \ddots \\ & & & & & -z & 1 & -z \\ & & & & & & -z & 1 \end{bmatrix} \cdot \begin{bmatrix} \psi_{h,-h}(z) \\ \psi_{h,-h+1}(z) \\ \vdots \\ \psi_{h,0}(z) \\ \vdots \\ \psi_{h,h}(z) \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}.$$

Compare [10], where a similar system is used.

Using Cramer's rule we get

$$(2.3) \quad \psi_{h,l}(z) = \frac{z^{|l|} a_{h-1}(z) a_{h-l-1}(z)}{a_{2h}(z)},$$

where  $a_i(z)$  denotes the determinant of the matrix in (2.2) with  $i+1$  rows. From [6], [9], [10] we know

$$(2.4) \quad a_i(z) = \frac{1}{1-v^2} \frac{1-v^{2i+4}}{(1+v^2)^{i+1}},$$

where the substitution  $z=v/(1+v^2)$  is used. Inserting (2.4) in (2.3) we get the desired result.  $\square$

Let  $\Psi_h(z) = \sum_{|l| \leq h} \psi_{h,l}(z)$  be the generating function of the number  $c_{n,h}$  of all paths with maximal deviation  $\leq h$ .

THEOREM 2.2.

$$(2.5) \quad \Psi_h(z) = \frac{(1+v^2)(1-v^{h+1})^2}{(1-v)^2(1+v^{2h+2})}.$$

PROOF.

$$\Psi_h(z) = \psi_{h,0}(z) + 2 \sum_{1 \leq l \leq h} \psi_{h,l}(z).$$

Using Theorem 2.1 we get the result by an elementary computation.  $\square$

THEOREM 2.3.

$$(2.6) \quad c_{n,h} = 2^n - 2 \sum_{\lambda \geq 0} \sum_{0 \leq l \leq h} \left[ \left[ \left\lfloor \frac{n-(h+2)}{2} \right\rfloor - 2\lambda(h+1) - l \right]^n + \left[ \left\lfloor \frac{n-(h+2)}{2} \right\rfloor - 2\lambda(h+1) - l \right]^n \right].$$

PROOF. Using Cauchy's integral formula we have

$$(2.7) \quad c_{n,h} = \frac{1}{2\pi i} \int_{(0+)}^{(0+)} \frac{dz}{z^{n+1}} \Psi_h(z) = \frac{1}{2\pi i} \int_{(0+)}^{(0+)} \frac{dv(1+v^2)^n(1+v)(1-v^{h+1})^2}{v^{n+1}(1-v)(1+v^{2h+2})}$$

where the substitution  $z = v/(1+v^2)$  was used. Hence  $c_{n,h}$  is the coefficient of  $v^n$  in

$$(2.8) \quad \frac{(1+v^2)^n(1+v)(1-v^{h+1})^2}{(1-v)(1+v^{2h+2})}.$$

Expanding the denominator we get

$$\frac{1}{(1-v)(1+v^{2h+2})} = \sum_{\lambda \geq 0} \sum_{l=0}^{2h+1} v^{4\lambda(h+1)+l}.$$

The result follows now by some elementary manipulations.  $\square$

As an alternative representation one obtains in a similar way:

THEOREM 2.4. *With the abbreviation*

$$B(n, k) = \sum_{l=0}^k \binom{n}{l},$$

$$(2.9) \quad c_{n,h} = 2 \sum_{\lambda \geq 0} (-1)^\lambda \left[ B\left(n, \left\lfloor \frac{n-(h+2)}{2} \right\rfloor - \lambda(h+1)\right) + B\left(n, \left\lfloor \frac{n-(h+2)}{2} \right\rfloor - \lambda(h+1)\right) \right]. \quad \square$$

A further representation for  $c_{n,h}$  can be obtained from (2.8) by partial fraction expansion.

THEOREM 2.5.

$$(2.10) \quad c_{n,h} = \frac{2^n}{h+1} \sum_{0 \leq l \leq h} (-1)^l \cos^n \left( \frac{2l+1}{2(h+1)} \pi \right) \operatorname{ctg} \left( \frac{2l+1}{4(h+1)} \pi \right). \quad \square$$

COROLLARY 2.6. For  $n \leq h$  we have

$$(2.11) \quad \sum_{0 \leq l \leq h} (-1)^l \cos^n \left( \frac{2l+1}{2(h+1)} \pi \right) \operatorname{ctg} \left( \frac{2l+1}{4(h+1)} \pi \right) = h+1. \quad \square$$

Since the total of paths of length  $n$  is  $2^n$ , we immediately obtain the following corollary:

COROLLARY 2.7. With  $d_{n,h}$  denoting the number of paths of length  $n$  and maximal deviation  $> h$ , we have

$$(2.12) \quad d_{n,h} = 2 \sum_{\lambda \geq 0} \sum_{0 \leq l \leq h} \left[ \left( \left\lfloor \frac{n-(h+2)}{2} \right\rfloor - 2\lambda(h+1) - l \right)^n + \left( \left\lfloor \frac{n-(h+2)}{2} \right\rfloor - 2\lambda(h+1) - l \right)^n \right]. \quad \square$$

Using Abel's summation formula, we get the following exact expression for  $D_n$ , the expected maximal deviation of a path of length  $n$ .

THEOREM 2.8.

$$(2.13) \quad D_n = 2^{-n} 2 \sum_{h \geq 0} \sum_{\lambda \geq 0} \sum_{0 \leq l \leq h} \left[ \left( \left\lfloor \frac{n-(h+2)}{2} \right\rfloor - 2\lambda(h+1) - l \right)^n + \left( \left\lfloor \frac{n-(h+2)}{2} \right\rfloor - 2\lambda(h+1) - l \right)^n \right]. \quad \square$$

The remainder of this section is devoted to the study of the asymptotic behaviour of  $D_n$ . For the sake of brevity we confine ourselves to the case of even  $n$ .

As a first step we have to approximate sums of binomial coefficients of the following kind ( $0 \leq a \leq b$ ):

$$(2.14) \quad S_{a,b} := 2^{-2n} \sum_{a \leq k \leq b} \binom{2n}{n+k}.$$

**THEOREM 2.9.** Let  $\varepsilon > 0$ . Assume  $0 \leq a \leq b = O(n^{1/2+\varepsilon})$  and  $k = O(n^{1/2+\varepsilon})$ . Furthermore let  $\operatorname{erfc}(x)$  be the well-known complement of the error function:

$$\operatorname{erfc}(x) = \frac{2}{\sqrt{\pi}} \int_x^{\infty} e^{-t^2} dt.$$

Then

$$(2.15) \quad \begin{aligned} \text{(i)} \quad S_{a,b} &= \frac{1}{2} \left[ \operatorname{erfc} \left( \frac{a - \frac{1}{2}}{\sqrt{n}} \right) - \operatorname{erfc} \left( \frac{b + \frac{1}{2}}{\sqrt{n}} \right) \right] (1 + O(n^{-1+\varepsilon})), \\ \text{(ii)} \quad S_{k,k} &= \left[ \operatorname{erfc} \left( \frac{k}{\sqrt{n}} \right) - \operatorname{erfc} \left( \frac{k + \frac{1}{2}}{\sqrt{n}} \right) + \frac{1}{2} \frac{ke^{-k^2/n}}{\sqrt{\pi n^3}} \right] (1 + O(n^{-1+\varepsilon})), \\ \text{(iii)} \quad S_{k,k} &= \left[ \operatorname{erfc} \left( \frac{k - \frac{1}{2}}{\sqrt{n}} \right) - \operatorname{erfc} \left( \frac{k}{\sqrt{n}} \right) - \frac{1}{2} \frac{ke^{-k^2/n}}{\sqrt{\pi n^3}} \right] (1 + O(n^{-1+\varepsilon})). \end{aligned}$$

**PROOF.** For the sake of brevity we only want to stress the main ideas of the proof. Following the presentation in [5, pp 179–182] one can distinguish three components of the error committed in the above approximations. The first component is effected by the approximation

$$2^{-2n} \binom{2n}{n+k} \sim 2^{-2n} \binom{2n}{n} e^{-k^2/n}.$$

The second component is due to the approximation of the middle term by Stirling's formula

$$2^{-2n} \binom{2n}{n} \sim \frac{1}{\sqrt{n\pi}}.$$

For this two error components it follows easily from the estimates given in [5] that we have

$$2^{-2n} \binom{2n}{n+k} = \frac{1}{\sqrt{n\pi}} e^{-k^2/n} (1 + O(n^{-1+\varepsilon})) \quad \text{for } k = O(n^{1/2+\varepsilon}).$$

The third component is caused by replacing summation over  $a \leq k \leq b$  by integration within the bounds  $a - \frac{1}{2}$  and  $b + \frac{1}{2}$ . Regarding part (i) of the above theorem it suffices to consider the approximation

$$\frac{1}{\sqrt{2\pi}} \sqrt{\frac{2}{n}} e^{-k^2/n} \sim \frac{1}{\sqrt{2\pi}} \int_{(k-\frac{1}{2})\sqrt{\frac{n}{2}}}^{(k+\frac{1}{2})\sqrt{\frac{n}{2}}} e^{-t^2/2} dt.$$

Estimating the difference between these two terms by appropriate chosen triangles



we are led to (i). As for part (ii) and (iii) a similar reasoning can be used to derive the indicated correction terms and to estimate the order of the error.  $\square$

The correction terms in the formulae (ii) and (iii) allow us to lower the order of the error from  $O(n^{-1/2+\varepsilon})$  to  $O(n^{-1+\varepsilon})$ .

**THEOREM 2.10.** *Let  $\sigma$  and  $\tau$  be the arithmetical functions defined by*

$$(2.16) \quad \sigma(m) = \sum_{\substack{m=(4\lambda+1)(h+1), \\ \lambda, h \geq 0}} 1, \quad \tau(m) = \sum_{\substack{m=(4\lambda+3)(h+1), \\ \lambda, h \geq 0}} 1.$$

Then

$$(2.17) \quad D_{2n} = 2 \sum_{m \geq 1} \left[ (\sigma(m) - \tau(m)) \operatorname{erfc} \left( \frac{m}{2\sqrt{n}} \right) \right] (1 + O(n^{-1+\varepsilon})) + \\ + \frac{1}{\sqrt{\pi n^3}} \sum_{m \geq 1} [(\sigma(m) - \tau(m)) m e^{-m^2/n}] (1 + O(n^{-1+\varepsilon})).$$

**PROOF.** Starting from (2.13), considering the cases  $h$  even or odd separately and regarding that  $\binom{2n}{n-k} = \binom{2n}{n+k}$ , we have

$$(2.18) \quad D_{2n} = 2^{-2n} 2 \sum_{\substack{h, \lambda \geq 0 \\ h \text{ even}}} \sum_{0 \leq l \leq h} 2 \left( n + \frac{h+2+4\lambda(h+1)+2l}{2} \right) + \\ + 2^{-2n} 2 \sum_{\substack{h, \lambda \geq 0 \\ h \text{ odd}}} \left[ \sum_{0 \leq l \leq h-1} 2 \left( n + \frac{h+3+4\lambda(h+1)+2l}{2} \right) + \right. \\ \left. + \left( n + \frac{h+1+4\lambda(h+1)}{2} \right) + \left( n + \frac{3h+3+4\lambda(h+1)}{2} \right) \right].$$

Now we approximate the last two terms in (2.8) by (ii) and (iii) of Theorem 2.9, respectively, and the remaining sums by (i) and obtain:

$$D_{2n} = 2 \sum_{h, \lambda \geq 0} \left[ \operatorname{erfc} \left( \frac{h+1+4\lambda(h+1)}{2\sqrt{n}} \right) - \operatorname{erfc} \left( \frac{3h+3+4\lambda(h+1)}{2\sqrt{n}} \right) \right] (1 + O(n^{-1+\varepsilon})) + \\ + \frac{2}{\sqrt{2\pi}} \sum_{k, \lambda \geq 0} \left[ \frac{k+1+4\lambda(k+1)}{\sqrt{2n^3}} \exp \left( \frac{-[k+1+4\lambda(k+1)]^2}{n} \right) - \right. \\ \left. - \frac{3k+3+4\lambda(k+1)}{\sqrt{2n^3}} \exp \left( \frac{-[3k+3+4\lambda(k+1)]^2}{n} \right) \right] (1 + O(n^{-1+\varepsilon})).$$

The error committed in the above approximation by extending the range of summation to infinity is exponentially small and therefore covered by the error terms.

The desired result now immediately follows by use of the arithmetical functions  $\sigma$  and  $\tau$ .  $\square$

The next lemma deals with the generating Dirichlet series of  $\sigma(m) - \tau(m)$ :

LEMMA 2.11.

$$(2.19) \quad \sum_{m \geq 1} \frac{\sigma(m) - \tau(m)}{m^z} = \frac{1}{4^z} \zeta(z) \left[ \zeta\left(z, \frac{1}{4}\right) - \zeta\left(z, \frac{3}{4}\right) \right],$$

where  $\zeta(z)$  is the zeta function of Riemann and  $\zeta(z, a)$  is the zeta function of Hurwitz (cf. e.g. [1], [12]).

PROOF. First we compute

$$\begin{aligned} \sum_{m \geq 1} \frac{\sigma(m)}{m^z} &= \sum_{m \geq 1} m^{-z} \sum_{\substack{m = (4\lambda+1)(h+1) \\ \lambda, h \geq 0}} 1 = \frac{1}{4^z} \sum_{\substack{m = (4\lambda+1)(h+1) \\ \lambda, h \geq 0}} \left(\lambda + \frac{1}{4}\right)^{-z} (h+1)^{-z} = \\ &= \frac{1}{4^z} \sum_{\lambda \geq 0} \left(\lambda + \frac{1}{4}\right)^{-z} \sum_{h \geq 0} (h+1)^{-z} = \frac{1}{4^z} \zeta\left(z, \frac{1}{4}\right) \zeta(z). \end{aligned}$$

The computation for  $\tau(m)$  is similar and yields

$$\sum_{m \geq 1} \frac{\tau(m)}{m^z} = \frac{1}{4^z} \zeta\left(z, \frac{3}{4}\right) \zeta(z). \quad \square$$

THEOREM 2.12. For all  $m > 0$  and  $n \rightarrow \infty$  we have

$$(2.20) \quad \sum_{m \geq 1} \left[ (\sigma(m) - \tau(m)) \operatorname{erfc}\left(\frac{m}{2\sqrt{n}}\right) \right] = \frac{1}{2} \sqrt{n\pi} - \frac{1}{4} + O(n^{-m}).$$

PROOF. By inversion of the Mellin transform of  $\operatorname{erfc}(x)$  we get

$$(2.21) \quad \operatorname{erfc}(x) = \frac{1}{2\pi i} \frac{1}{\sqrt{\pi}} \int_{c-i\infty}^{c+i\infty} \frac{1}{z} \Gamma\left(\frac{z+1}{2}\right) x^{-z} dz, \quad x > 0, \quad c > 0,$$

which can be found in [4, p. 325]. Thus (2.20) equals

$$\sum_{m \geq 1} (\sigma(m) - \tau(m)) \frac{1}{2\pi i} \frac{1}{\sqrt{\pi}} \int_{c-i\infty}^{c+i\infty} \frac{1}{z} \Gamma\left(\frac{z+1}{2}\right) \left(\frac{m}{2\sqrt{n}}\right)^{-z} dz.$$

The sum may be placed inside the integral, since convergence is absolutely well behaved (cf. [8, p. 133]) and thus the last expression equals

$$(2.22) \quad \frac{1}{2\pi i} \frac{1}{\sqrt{\pi}} \int_{c-i\infty}^{c+i\infty} \frac{1}{z} \left(\frac{\sqrt{n}}{2}\right)^z \zeta(z) \left[ \zeta\left(z, \frac{1}{4}\right) - \zeta\left(z, \frac{3}{4}\right) \right] \Gamma\left(\frac{z+1}{2}\right) dz.$$

By a well-known method it can be shown that the line of integration can be shifted to the left as far as we please if we only take the residues into account, yielding an error of  $O(n^{-m})$  for all  $m > 0$ . It is important to notice that, by cancellation of the

poles of Hurwitz' zeta functions,  $\zeta\left(z, \frac{1}{4}\right) - \zeta\left(z, \frac{3}{4}\right)$  is an entire function! At  $z=1$  there is a simple pole with residue

$$(2.23) \quad \frac{\sqrt{n\pi}}{2}.$$

At  $z=0$  there is a simple pole with residue

$$(2.24) \quad -\frac{1}{4}.$$

At  $z=-(2k+1)$ ,  $k \in \mathbb{N}_0$ , the residues are zero, because

$$(2.25) \quad \begin{aligned} & \zeta\left(-(2k+1), \frac{1}{4}\right) - \zeta\left(-(2k+1), \frac{3}{4}\right) = \\ & = \frac{1}{2(k+1)} \left[ -B_{2(k+2)}\left(\frac{1}{4}\right) + B_{2(k+2)}\left(1 - \frac{1}{4}\right) \right] = 0 \end{aligned}$$

where  $B_m(x)$  is the  $m$ -th Bernoulli polynomial (cf. [3, p. 49]). Summing up we get the desired result.  $\square$

THEOREM 2.13. For  $m > 0$  and  $n \rightarrow \infty$ ,

$$(2.26) \quad \sum_{m \geq 1} (\sigma(m) - \tau(m)) m e^{-m^2/n} = \frac{n\pi}{8} + O(n^{-m}).$$

PROOF. Using the well-known formula

$$(2.27) \quad e^{-x} = \frac{1}{2\pi i} \int_{c-i\infty}^{c+i\infty} \Gamma(z) x^{-z} dz, \quad x > 0, \quad c > 0$$

a similar method as in the proof of Theorem 2.12 (cf. e.g. [2]) yields the result.  $\square$

We want to summarize the results of this section in the following theorem.

THEOREM 2.14. The expected maximal deviation of a path of length  $n$ ,  $n$  even, is given by

$$D_n = \sqrt{\frac{n\pi}{2}} - \frac{1}{2} + O(n^{-1/2+\varepsilon}) \quad \text{for all } \varepsilon > 0 \quad \text{and } n \rightarrow \infty.$$

PROOF. The application of Theorems 2.12 and 2.13 to (2.17) yields the result.  $\square$

### 3. The expected maximal span

Let  $\Psi_h(z)$  be the generating function whose  $n$ -th coefficient gives the number of paths of length  $n$  with maximal span  $\leq h$ . Such a path is shown in the following figure.

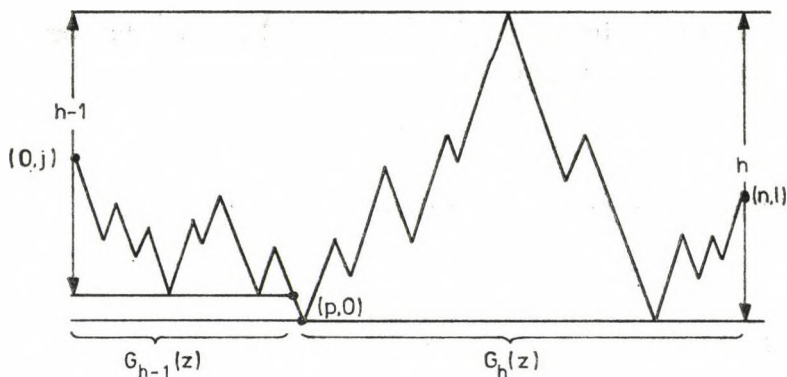


Fig. 2

To obtain an explicit expression for  $\Psi_h(z)$  we proceed as follows: The minimal value of the ordinate of a path enumerated by  $\Psi_h(z)$  is interpreted as the 0-level. Considering the point  $P=(p, 0)$ , where the minimal level is reached for the first time (in a left-to-right-sense), the desired generating function can be obtained as the convolution of the two generating functions describing the left and the right part: For the right part we have to count the number of nonnegative paths with height  $\leq h$ . For the left part we proceed from  $P$  to the left. If  $p=0$ , the contribution is 1; otherwise the first step leads to the point  $(p-1, 1)$ . For the number of the remaining paths we have to count the nonnegative paths of height  $\leq h-1$ .

Using the generating function  $G_h(z)$  of nonnegative paths with height  $\leq h$ , defined in [6], we have

THEOREM 3.1.

$$(3.1) \quad \Psi_h(z) = (1 + zG_{h-1}(z))G_h(z). \quad \square$$

The generating functions  $G_h(z)$  are given by

$$(3.2) \quad G_h(z) = \frac{(1+v^2)(1-v^{h+1})}{(1-v)(1+v^{h+2})}$$

where the substitution  $z=v/(1+v^2)$  was used. We would like to emphasize that (3.1) holds also for  $h=0$ , since  $G_{-1}(z)=0$ . Substituting (3.2) in (3.1) we get

$$(3.3) \quad \Psi_h(z) = \frac{(1+v^2)(1-v^{h+1})(1-v^{h+2})}{(1-v)^2(1+v^{h+1})(1+v^{h+2})}, \quad h \geq 0.$$

Let  $d_{n,h}$  denote the number of paths of length  $n$  and maximal span  $>h$  and  $s_n = \sum_{h=0}^n d_{n,h}$ .

THEOREM 3.2.  $s_n$  is the coefficient of  $v^n$  in

$$(3.4) \quad 2 \frac{(1+v^2)^n(1+v)v}{(1-v)^2}.$$

PROOF. The generating function  $\chi_h(z)$  of all paths with maximal span  $> h$  is the difference between the generating function of all paths (i.e.  $\frac{1}{1-2z} = \frac{1+v^2}{(1-v)^2}$ ) and  $\Psi_h(z)$ . Hence

$$(3.5) \quad \chi_h(z) = 2 \frac{1+v^2}{(1-v)^2} \frac{(1+v)v^{h+1}}{(1+v^{h+1})(1+v^{h+2})}.$$

By Cauchy's integral formula we obtain

$$\begin{aligned} s_n &= \sum_{h=0}^n d_{n,h} = \sum_{h=0}^n \frac{1}{2\pi i} \int_{z^{n+1}}^{(0,+)} \frac{dz}{z^{n+1}} \chi_h(z) = \\ &= \frac{1}{2\pi i} \int_{v^{n+1}}^{(0,+)} \frac{dv}{v^{n+1}} 2 \frac{(1+v)^2(1+v^2)^n}{(1-v)^2} \sum_{h=0}^n \left( \frac{1}{1+v^{h+2}} - \frac{1}{1+v^{h+1}} \right) = \\ &= \frac{1}{2\pi i} \int_{v^{n+1}}^{(0,+)} \frac{dv}{v^{n+1}} 2 \frac{(1+v)^2(1+v^2)^n}{(1-v)^2} \frac{v}{1+v}, \end{aligned}$$

which immediately leads to the desired result.  $\square$

THEOREM 3.3.

$$s_{2n} = 2 \left[ (4n+1) \binom{2n-1}{n} - 2^{2n-1} \right],$$

$$s_{2n+1} = 2 \left[ 4(2n+1) \binom{2n-1}{n} - 2^{2n} \right].$$

PROOF. This can be obtained by (3.4) using the following identities (cf. [11, p. 34]):

$$\sum_{k=0}^n \binom{2n}{n-k} = 2^{2n-1} + \binom{2n-1}{n}, \quad \sum_{k=0}^n \binom{2n+1}{n-k} = 2^{2n},$$

$$\sum_{k=0}^n k \binom{2n}{n-k} = n \binom{2n-1}{n}, \quad \sum_{k=0}^n k \binom{2n+1}{n-k} = (2n+1) \binom{2n-1}{n} - 2^{2n-1}. \quad \square$$

THEOREM 3.4. The expected maximal span of a path of length  $n$  is given by

$$\sqrt{\frac{8n}{\pi}} - 1 + O(n^{-1/2}), \quad n \rightarrow \infty.$$

PROOF. Since the expected maximal span is  $2^{-n}s_n$ , we obtain the result by Theorem 3.3 and Stirling's formula.  $\square$

#### 4. The expected onesided height

Let  $\Psi_{h,k;l}(z)$  be the generating function whose  $n$ -th coefficient corresponds to all paths with  $-k \leq a_i \leq h$  for  $0 \leq i \leq n$ ,  $a_0 = 0$  and  $a_n = l$ ;  $\Psi_h(z)$  corresponds to all paths with  $a_i \leq h$  and  $a_n = 0$ .

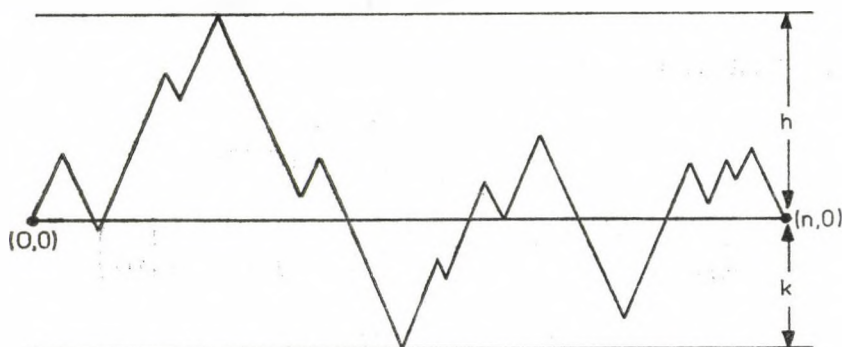


Fig. 3

THEOREM 4.1. With  $z = v/(1+v^2)$ ,

$$(4.1) \quad \Psi_h(z) = \frac{1+v^2}{1-v^2} (1-v^{2h+2}).$$

PROOF. By a similar argument as in the proof of Theorem 2.1, we find that

$$(4.2) \quad \Psi_{h,k;0}(z) = \frac{a_{h-1}(z)a_{k-1}(z)}{a_{h+k}(z)}$$

with the determinants  $a_l(z)$  of Theorem 2.1. Hence

$$(4.3) \quad \Psi_{h,k;0}(z) = \frac{1+v^2}{1-v^2} [1-v^{2h+2}] \frac{1-v^{2k+2}}{1-v^{2h+2k+4}}.$$

Since  $\Psi_h(z) = \lim_{k \rightarrow \infty} \Psi_{h,k;0}(z)$  we have the result for small values of  $v$  (and thus small values of  $z$ ) and by continuation for all values of  $z$  in the circle of convergence of  $\Psi_h(z)$ .  $\square$

Let  $s_n$  denote the sum of heights of all paths of length  $2n$ .

THEOREM 4.2.

$$(4.4) \quad s_n = 2^{2n-1} - \binom{2n-1}{n}.$$

PROOF. Let  $D(z) = (1-4z^2)^{-1/2} = \sum_{n \geq 0} \binom{2n}{n} z^{2n}$  be the generating function of



all paths of length  $2n$ . Then, as usual,  $s_n$  is the coefficient of  $z^{2n}$  in

$$\begin{aligned}
 (4.5) \quad & \sum_{h \equiv 0} (D(z) - \Psi_h(z)) = \\
 & = \sum_{h \equiv 0} \left( \frac{1+v^2}{1-v^2} - \frac{1+v^2}{1-v^2} (1-v^{2h+2}) \right) = \frac{1+v^2}{1-v^2} \sum_{h \equiv 0} v^{2h+2} = v^2 \frac{1+v^2}{(1-v^2)^2}. \\
 & s_n = \frac{1}{2\pi i} \int_{(0+)} \frac{dz}{z^{2n+1}} v^2 \frac{1+v^2}{(1-v^2)^2} = \frac{1}{2\pi i} \int_{(0+)} \frac{dv}{v^{2n+1}} \frac{v^2(1+v^2)^{2n}}{1-v^2}.
 \end{aligned}$$

Hence  $s_n$  is the coefficient of  $u^n$  in

$$\frac{u(1+u)^{2n}}{1-u},$$

which is

$$(4.6) \quad \sum_{\lambda \equiv 1} \binom{2n}{n-\lambda}. \quad \square$$

COROLLARY 4.3. *The expected onesided height of a path of length  $2n$  is given by*

$$(4.7) \quad \frac{1}{2} \sqrt{\pi n} - \frac{1}{2} + O(n^{-1/2}) \quad \text{for } n \rightarrow \infty.$$

PROOF. Apply Stirling's approximation formula to  $\binom{2n}{n}^{-1} s_n$ .  $\square$

## 5. Concluding remarks

We would like to mention that Theorem 2.14 also holds true for  $n$  odd. This could be shown by additionally considering the case  $n$  odd in the derivation of Theorem 2.14. But it suffices for our purposes to observe that the expected maximal deviation is a strictly increasing function of path length. Now — since the difference  $D_{2n+2} - D_{2n}$  is of order  $O(n^{-1/2})$  — the validity of Theorem 2.14 for all  $n$  immediately follows.

Distinguishing the three components of the approximation error dealt with in the proof of Theorem 2.9, a good balance of their order was achieved by taking the correction terms into account in Theorem 2.9 (ii) and (iii). By that means we were able to lower the order of the error term in Theorem 2.14 from  $O(n^\epsilon)$  to  $O(n^{-1/2+\epsilon})$  and thereby the absolute term could be preserved. In principle better approximations could be obtained by use of Euler's summation formula. However, it can be seen from the following table that the asymptotic formulae derived in Section 2 show an accuracy meeting most practical requirements even for small  $n$ .

Furthermore, a point of some methodical interest should be emphasized: Though the three problems treated in Sections 2 to 4 seem to be very akin, in the study of their asymptotic behaviour different methods had to be used.

In Section 3 and 4 the exact formulae were reduced to a relatively simple form by the use of some combinatorial identities. To derive the corresponding asymptotic formulae essentially Stirling's formula had to be applied.

$\begin{smallmatrix} D_n \\ n \end{smallmatrix}$	exactly (Theorem 2.8)	asymptotically (Theorem 2.14)
10	3.53	3.46
20	5.15	5.11
30	6.40	6.36
40	7.46	7.43
50	8.39	8.36
60	9.24	9.21
70	10.01	9.99
80	10.73	10.71
90	11.41	11.39
100	12.05	12.03

For the study of the asymptotic behaviour of the expected maximal deviation treated in Section 2, the so-called  *$\Gamma$ -function method*, used e.g. in [2], had to be modified: instead of the exponential function the complement of the error function was used. In addition to Riemann's zeta-function also Hurwitz' zeta-function came into play.

#### REFERENCES

- [1] APOSTOL, T. M., *Introduction to Analytic Number Theory*, Springer, New York—Heidelberg—Berlin, 1976. *MR* 55 # 7892.
- [2] DE BRUIJN, N. G., KNUTH, D. E. and RICE, S. O., The average height of planted plane trees, in: *Graph theory and computing*, R. C. Read, Ed., 15—22. Academic Press, New York—London, 1972. *MR* 48 # 8280.
- [3] COMTET, L., *Advanced Combinatorics*, Cloth Edition, D. Reidel, Dordrecht, 1974. *MR* 57 # 124.
- [4] ERDÉLYI, A., MAGNUS, W., OBERHETTINGER, F. and TRICOMI, F. G., *Tables of Integral Transform*, Vol. 1, McGraw-Hill, New York—Toronto—London, 1954. *MR* 15—868.
- [5] FELLER, W., *An introduction to probability theory and its applications*, Vol. 1, 3rd ed., J. Wiley, New York, 1968. *MR* 37 # 3604.
- [6] KEMP, R., On the average depth of a prefix of the Dycklanguage  $D_1$ , *Discrete Math.* 36 (1981), 155—170.
- [7] KIRSCHENHOFER, P. and PRODINGER, H., On the average hyperoscillations of planted plane trees, *Combinatorica* 2 (1982), 177—186.
- [8] KNUTH, D. E., *The Art of Computer Programming*, Vol. 3, Addison-Wesley, Reading, Massachusetts, 1973. *MR* 56 # 4281.
- [9] PANNY, W., The expected depth of a pushdown store with three operations and symmetrical probabilities, *J. Comb. Inf. Syst. Sci.* 7 (1982), 38—47.
- [10] PRODINGER, H., The average height of a stack where three operations are allowed and some related problems, *J. Combin. Inform. System Sci.* 5 (1980), 287—304. *MR* 82 m: 05013.
- [11] RIORDAN, J., *Combinatorial Identities*, J. Wiley, New York, 1968. *MR* 38 # 53.
- [12] WHITTAKER, E. T. and WATSON, G. N., *A course of modern analysis*, 4th ed., Cambridge, University Press, 1927.

(Received January 28, 1982)

INSTITUT FÜR STATISTIK  
WIRTSCHAFTSUNIVERSITÄT WIEN  
AUGASSE 2—6  
A—1090 WIEN

INSTITUT FÜR ALGEBRA UND  
DISKRETE MATHEMATIK  
TECHNISCHE UNIVERSITÄT  
WIEDNER HAUPTSTRASSE 8—10  
A—1040 WIEN  
AUSTRIA

## ESTIMATION OF THE PARAMETERS OF BURR DISTRIBUTION BASED ON ORDER STATISTICS

MUNIR AHMAD

### Abstract

Asymptotically best unbiased estimators of the two Burr parameters  $\alpha$  and  $\beta$  based on some selected order statistics from a sample of size  $n$  are considered. It is shown that the minimum variance of the estimator of  $\beta$  for the known  $\alpha$  is  $1.54416\beta^2/n$  with efficiency of about 64% as compared to the maximum likelihood estimator of  $\beta$ . The minimum variance of the estimator of  $\alpha$  for known  $\beta$  is attained when 76th sample order statistics is used. When  $\alpha$  and  $\beta$  are unknown, variances are essentially minimized using three order statistics (.24, .54, .77) for  $\alpha$  and (.27, .61, .81) for  $\beta$ . An example is given to illustrate the method.

### 1. Introduction

Consider a random sample of size  $n$  from the Burr distribution first given by Burr (1942) having the probability density function

$$(1) \quad f(x) = \alpha\beta x^{\alpha-1}/(1+x^\alpha)^{\beta+1}, \quad x > 0, \quad \beta \geq 1$$

and the distribution function

$$(2) \quad F(x) = 1 - (1+x^\alpha)^{-\beta}, \quad x > 0.$$

Properties of the distribution (1) have been investigated by Austin (1971), Burr (1968), Burr and Cislak (1968) and Hatke (1949). Recently Austin (1971) has applied the distribution to the control charts for the maximum and minimum value in sampling from a normal distribution. In modelling tensile loads on a machine element, Weibull and log normal distributions have been tried (see Bury, 1975). In this paper we have used some order statistics to estimate the Burr parameters.

### 2. Estimation

Let  $x_1 \leq x_2 \leq \dots \leq x_n$  be the sample ordered statistics and let  $u_1, u_2, \dots, u_n$  be the corresponding population quantities. Let  $x_{n_1} < x_{n_2} < \dots < x_{n_k}$ ,  $k=1, 2, \dots$  be the selected  $k$  sample quantities where  $k < n$ . A set of  $k$  real numbers such that  $0 = \lambda_0 < \lambda_1 < \dots < \lambda_{k+1} = 1$  is called a spacing. We define  $\lambda_i$  by  $\lambda_i = F(u_i)$ ,  $i=1, 2, \dots, k$ . Then  $u_i = F^{-1}(\lambda_i)$  and therefore from (2)

$$(3) \quad u_i = [(1 - \lambda_i)^{-1/\beta} - 1]^{1/\alpha}$$

and from (1)

$$(4) \quad f[F^{-1}(\lambda_i)] = \alpha\beta(1-\lambda_i)^{1+1/\beta} u_i^{\alpha-1}.$$

Taking the logarithm of (3) and (4), we get

$$(5) \quad \ln u_i = 1/\alpha \ln [(1-\lambda_i)^{-1/\beta} - 1]$$

and

$$(6) \quad \ln f[F^{-1}(\lambda_i)] = \ln(\alpha\beta) + (1+1/\beta) \ln(1-\lambda_i) + (\alpha-1) \ln u_i.$$

Using some selected order statistics, we shall obtain best unbiased estimators of  $\alpha$  and  $\beta$  by minimizing the asymptotic variances of the estimators.

*Case 1.* ( $\beta$  is known): Given a spacing  $\{\lambda_i\}$ , the estimator of  $\alpha$  when  $\beta = \beta_0$ , using (3) is

$$(7) \quad \tilde{\alpha}_k = (k)^{-1} \sum_{i=1}^k \{\ln [(1-\lambda_i)^{-1/\beta_0} - 1] y_i\} = (k)^{-1} \sum_{i=1}^k a_i y_i$$

where  $y_i = (\ln x_{n_i})^{-1}$ ,  $a_i = \ln [(1-\lambda_i)^{-1/\beta_0} - 1]$  and  $x_{n_i}$  is a selected sample ordered statistic corresponding to  $u_{n_i}$ . The asymptotic variance of  $\tilde{\alpha}_k$  is given by

$$(8) \quad \text{var}(\tilde{\alpha}_k) = k^{-2} \sum a_i^2 \text{var} y_i + \sum_{i \neq j} a_i a_j \text{cov}(y_i, y_j),$$

where

$$\text{var}(y_i) = y_i^4 x_{n_i}^{-2} \text{var}(x_{n_i})$$

and

$$\text{cov}(y_i, y_j) = y_i^2 y_j^2 (x_{n_i} x_{n_j})^{-1} \text{cov}(x_{n_i}, x_{n_j}).$$

It is well-known [see Mosteller (1946)] that given a spacing, the joint distribution of the  $k$  sample quantities is asymptotically normal. Suppose  $k=i$ . The  $i$ th sample quantity is asymptotically normal with mean  $u_i$  and variance  $\lambda_i(1-\lambda_i)/[nf^2(x_{n_i})]$ . Minimizing  $\text{var}(\tilde{\alpha}_1)$  with respect to  $\lambda_i$ , we get

$$(9) \quad (1-2\lambda_i)[1-(1-\lambda_i)^{1/\beta_0}][\ln(1-\lambda_i)^{-1/\beta_0} - 1] + \lambda_i/\beta_0 = 0.$$

Table 1 gives the solution of  $\lambda_i$  of equation (9) that minimizes  $\text{var}(\tilde{\alpha}_1)$  for some given values of  $\beta$ .

Table 1

$\beta$	$\lambda_i$	$a_i$
1.0	.823958	1.54340
1.5	.869111	1.05748
2.0	.904915	0.80773
2.5	.931742	0.65567
3.0	.951335	0.55332
3.5	.965447	0.47972
4.0	.975530	0.42421
5.0	.987780	0.34583
6.0	.993912	0.29282
7.0	.996969	0.25439
8.0	.998491	0.22514
9.0	.999248	0.20192
10.0	.999626	0.18355

In case of large  $\beta$ , the largest observation may be used to estimate  $\alpha$ . For  $k=2, 3, \dots, \lambda_i$  can be determined by minimizing (8) numerically.

Case 2. ( $\alpha$  is known): Suppose  $\alpha = \alpha_0$ . For a given spacing  $\{\lambda_i\}$  satisfying (2), we find that

$$(10) \quad \tilde{\beta}_k = k^{-1} \sum_{i=1}^k b_i z_i$$

where

$$b_i = \ln(1 - \lambda_i), \quad \text{and} \quad z_i = -\{\ln(1 + x_{n_i}^{\alpha_0})\}^{-1}.$$

Its asymptotic variance is given by

$$\text{var}(\tilde{\beta}_k) = k^{-2} \left[ \sum_i b_i^2 \text{var}(z_i) + \sum_{i \neq j} b_i b_j \text{cov}(z_i, z_j) \right],$$

where

$$\text{var}(z_i) = z_i^4 f_i^2 [\beta(1 - \lambda_i)]^{-2} \text{var}(x_{n_i}), \quad f_i = f(x_{n_i})$$

and

$$\text{cov}(z_i, z_j) = z_i^2 z_j^2 f_i f_j \beta^{-2} [(1 - \lambda_i)(1 - \lambda_j)]^{-1} \text{cov}(x_{n_i}, x_{n_j}).$$

Let  $k=i$ ,  $\tilde{\beta}_i = b_i z_i$  and the asymptotic variance is given by

$$(11) \quad \text{var}(\tilde{\beta}_i) = \left[ \frac{\ln(1 - \lambda_i)}{\ln^2(1 + x_{n_i}^{\alpha_0})} \frac{\alpha_0 x_{n_i}^{\alpha_0 - 1}}{1 + x_{n_i}^{\alpha_0}} \right]^2 \text{var}(x_i).$$

Using  $\text{var}(x_i) = \lambda_i(1 - \lambda_i)/nf_i^2$ ,

$$f_i = \alpha_0 \beta x_{n_i}^{\alpha_0 - 1} (1 - \lambda_i) / [1 + x_{n_i}^{\alpha_0}] \quad \text{and} \quad \ln(1 + x^{\alpha_0}) = -(1/\beta) \ln(1 - \lambda_i).$$

We have

$$(12) \quad \text{var}(\tilde{\beta}_i) = \beta^2 \lambda_i / [n(1 - \lambda_i) \ln^2(1 - \lambda_i)].$$

We choose  $\lambda_i$  by minimizing  $\text{var}(\tilde{\beta}_i)$  with respect to  $\lambda_i$  and derive an equation in  $\lambda_i$  only

$$2\lambda_i + \ln(1 - \lambda_i) = 0$$

giving a value of  $\lambda_i = 0.79681$ .

The minimum variance of  $\tilde{\beta}$  for any  $\alpha_0$  is  $1.54416 \beta^2/n$ . The maximum likelihood estimate of  $\beta$  for  $\alpha = \alpha_0$  from a sample of size  $n$  is

$$\hat{\beta} = n / \sum_{i=1}^n \ln(1 + x_i^{\alpha_0})$$

and its asymptotic variance is  $\beta^2/n$ . The efficiency of  $\tilde{\beta}$  as compared to  $\hat{\beta}$  is given by

$$E = (1 - \lambda_i) \ln(1 - \lambda_i) / \lambda_i$$

and is about 64% for  $\lambda_i = 0.79681$ .

For  $k=2, 3, \dots, \lambda_i$  can be determined such that (11) is minimized. Some special cases have been investigated separately.

**Case 3.** ( $\alpha$  and  $\beta$  are unknown): Suppose  $\alpha$  and  $\beta$  are unknown. We write  $y_{n_i} = \ln x_{n_i} = (1/\alpha)[(1-\lambda_i)^{-1/\beta} - 1]$  and

$$\ln f_{n_i} = \ln \tilde{\alpha} + \ln \tilde{\beta} + (\alpha - 1)y_{n_i} + (1 + 1/\tilde{\beta}) \ln(1 - \lambda_i), \quad i = 1, 2, 3.$$

Using  $\ln f_{n_i}$ ,  $i = 1, 2, 3$ , we have

$$l_1 = (\tilde{\alpha} - 1)y_1 + (1 + 1/\tilde{\beta})\alpha_1$$

$$l_2 = (\tilde{\alpha} - 1)y_2 + (1 + 1/\tilde{\beta})\alpha_2$$

where

$$l_1 = \ln(f_{n_1}/f_{n_2}), \quad l_2 = \ln(f_{n_2}/f_{n_3})$$

$$\alpha_1 = \ln[(1 - \lambda_1)/(1 - \lambda_2)], \quad \alpha_2 = \ln[(1 - \lambda_2)/(1 - \lambda_3)]$$

$$y_1 = y_{n_1} - y_{n_2}, \quad y_2 = y_{n_2} - y_{n_3} \quad \text{and} \quad \lambda_i = F(x_{n_i}).$$

Solving for  $\tilde{\alpha}$  and  $\tilde{\beta}$ , we get the estimator of  $(\alpha, \beta)$  as

$$(13) \quad \tilde{\alpha} = 1 + (l_1\alpha_2 - l_2\alpha_1)/(y_1\alpha_2 - y_2\alpha_1)$$

and

$$(14) \quad \tilde{\beta} = (\alpha_1 y_2 - \alpha_2 y_1)/[(l_1 - \alpha_1)y_2 - (l_2 - \alpha_2)y_1].$$

The asymptotic variances of  $\tilde{\alpha}$  and  $\tilde{\beta}$  are

$$(15) \quad \text{var}(\tilde{\alpha}) = (\alpha - 1)^2 \sum_{i=1}^4 \sum_{j=1}^4 a_i a_j C(\eta_i, \eta_j)$$

and

$$(16) \quad \text{var}(\tilde{\beta}) = \beta^2 \sum_{i=1}^4 \sum_{j=1}^4 b_i b_j C(\eta_i, \eta_j)$$

where

$$C(\eta_i, \eta_j) = \begin{cases} \text{Var}(\eta_i), & j = i \\ \text{Cov}(\eta_i, \eta_j), & j \neq i \end{cases}$$

$$\eta_1 = l_1, \quad \eta_2 = l_2, \quad \eta_3 = y_1, \quad \eta_4 = y_2, \quad a_1 = \alpha_2/(l_1\alpha_2 - l_2\alpha_1)$$

$$a_2 = -\alpha_1/(l_1\alpha_2 - l_2\alpha_1)$$

$$a_3 = -\alpha_2/(y_1\alpha_2 - y_2\alpha_1)$$

$$a_4 = \alpha_1/(y_1\alpha_2 - y_2\alpha_1)$$

$$b_1 = -y_2/[(l_1 - \alpha_1)y_2 - (l_2 - \alpha_2)y_1]$$

$$b_2 = y_1/[(l_1 - \alpha_1)y_2 - (l_2 - \alpha_2)y_1]$$

$$b_3 = (l_2 - \alpha_2)/[(l_1 - \alpha_1)y_2 - (l_2 - \alpha_2)y_1] - a_3$$

$$b_4 = -a_4 - (l_1 - \alpha_1)/[(l_1 - \alpha_1)y_2 - (l_2 - \alpha_2)y_1].$$



Minimum variance can be obtained only numerically using variance values of unknown quantities involved in the expressions (15) and (16). For various values of  $\{\lambda_i\}$  the variance functions at (15) and (16) were computed on the IBM 370 computer. For  $k=3$ , we find that the variance functions are essentially minimized by choosing (0.24, 0.54, 0.77) for  $\text{var}(\hat{\alpha})$  and by choosing (0.27, 0.61, 0.81) for  $\text{var}(\hat{\beta})$ .

### 3. Example

To illustrate the method, suppose that the machine element is subjected to 100 random tensile loads during its mission  $x$  and that the load has been represented by the Burr model with  $\alpha=6$  and  $\beta=4$ . If  $\beta$  is known to have a value of 4, then  $\alpha$  is estimated from (7) when  $k=i$  and  $\lambda_i=0.97553$  (see Table 1). The values of  $x_{n_i}$  and  $y_i$  are  $x_{.98}=1.075$ ,  $y_i=0.93023$  and  $a_i=0.42421$ .  $\hat{\alpha}=5.87$  and  $\text{var}(\hat{\alpha})=13.07$ . If  $\alpha$  is known to have a value of 6, then  $\beta$  is estimated from (10) when  $k=1$ ,  $\lambda_i=0.79681$ . The values of  $x_{n_i}=0.853$ ,  $z_i=-3.0689$  and  $b_i=-1.5936$ . Then  $\hat{\beta}=4.89$  and  $\text{var}(\hat{\beta})=0.3696$ . In this case, the maximum likelihood estimation of  $\beta$  is  $\hat{\beta}=4.23$  and  $\text{var}(\hat{\beta})=0.1786$ . If  $\beta=4$ , then efficiency of  $\hat{\beta}$  is about 64%. If we use  $\hat{\beta}$  and  $\hat{\beta}$  in the variance formulae, respectively, the efficiency is about 48.3%. If  $\alpha$  and  $\beta$  are unknown then using (13) and (14) as estimating equations we have for  $\alpha$ , the values of  $x_{n_1}$ ,  $x_{n_2}$  and  $x_{n_3}$  are  $x_{n_1}=.585$ ,  $x_{n_2}=.744$  and  $x_{n_3}=.840$ . Then  $\hat{\alpha}=5.53$  and  $\text{var}(\hat{\beta})=20.31$  and for  $\beta$ , the values of  $x_{n_1}$ ,  $x_{n_2}$  and  $x_{n_3}$  are  $x_{n_1}=.612$ ,  $x_{n_2}=.772$  and  $x_{n_3}=.878$ . Then  $\hat{\beta}=3.42$  and  $\text{var}(\hat{\beta})=.2511$ .

### REFERENCES

- [1] AUSTIN, J. A. JR., Control chart constants for largest and smallest in sampling from a normal distribution using the generalized Burr distribution, *Technometrics* **15** (1971), 931—933.
- [2] BURR, I. W., Cumulative frequency functions, *Ann. Math. Statistics* **13** (1942), 215—232. *MR* **4**—99.
- [3] BURR, I. W., On a general system of distribution III. The sample range, *J. Amer. Statist. Assoc.* **63** (1968), 636—643. *MR* **38** #2860.
- [4] BURR, I. W. and CISLAK, P. J., On a general system of distributions. Its curve-shape characteristics — II The sample median, *J. Amer. Statist. Assoc.* **63** (1968), 627—635. *MR* **38** #2859.
- [5] BURY, K. V., *Statistical Models in Applied Science*, John Wiley & Sons, New York, 1975. *MR* **55** #4574.
- [6] HATKE, M. A., A certain cumulative probability function, *Ann. Math. Statistics* **20** (1949), 461—463. *MR* **11**—41.
- [7] MOSTELLER, F., On some useful "inefficient" statistics, *Ann. Math. Statistics* **17** (1946), 377—408. *MR* **8**—477.

(Received April 14, 1982)

DEPARTMENT OF MATHEMATICS  
UNIVERSITY OF PETROLEUM AND MINERALS  
P.O. BOX 476  
DHAHRAN  
SAUDI ARABIA



# INDECOMPOSABLE TRIPLE SYSTEMS WITH $\lambda=4$

CHARLES J. COLBOURN<sup>1</sup> and ALEXANDER ROSA<sup>2</sup>

## Abstract

A triple system  $B[3, \lambda; v]$  is indecomposable if it is not the union of two triple systems  $B[3, \lambda'; v]$  and  $B[3, \lambda''; v]$  with  $\lambda = \lambda' + \lambda''$ . Recursive constructions, along with direct constructions for eight small orders, are used to establish that indecomposable triple systems with  $\lambda=4$  exist if and only if  $v \equiv 0, 1 \pmod{3}$ ,  $v \geq 10$ . These triple systems have no repeated blocks.

## 1. Introduction

A *balanced incomplete block design*, denoted  $B[k, \lambda; v]$ , is a pair  $(V, B)$ ;  $V$  is a  $v$ -set of *elements* and  $B$  is a collection of  $k$ -element subsets of  $V$  called *blocks*. Each 2-subset of  $V$  appears in precisely  $\lambda$  blocks. When  $k=3$ , block designs are called *triple systems*. When  $B$  contains no repeated blocks (i.e., contains no block twice), we use the notation  $NB[k, \lambda; v]$ .

One common construction technique for block designs is to combine small block designs to form larger ones. A trivial way to do this combines a  $B[k, \lambda; v]$  design  $(V, B)$  with a  $B[k, \lambda'; v]$  design  $(V, B')$ ; their union  $(V, B \cup B')$  is a  $B[k, \lambda + \lambda'; v]$  design. Many nonisomorphic designs can be formed in this way; such designs are called *decomposable*. All designs can be constructed easily from indecomposable designs, and thus it is reasonable to ask for which  $v, k$ , and  $\lambda$  an indecomposable design exists. Informally, this is asking what the building blocks are for designs with large  $\lambda$ .

Kramer [7] initiated research in this area, showing that indecomposable  $NB[3, 2; v]$  designs exist for all  $v \equiv 0, 1 \pmod{3}$ ,  $v \geq 6$ , except for  $v=7$ , and indecomposable  $NB[3, 3; v]$  designs exist for all  $v \equiv 1 \pmod{2}$ ,  $v \geq 5$ . Since that time, Billington [1, 2] has studied techniques applicable to larger  $k$  when repeated blocks are allowed; as she notes [2], these techniques apply when  $\lambda < k$ .

Recently, Colbourn and Colbourn [3] produced indecomposable  $B[3, \lambda; v]$  designs for all odd  $\lambda$ , as a byproduct of an *NP*-completeness result. This construction can also be adapted to produce an infinite family of indecomposable  $NB[3, 4; v]$  designs [4].

In this paper, we study the case  $k=3$  and  $\lambda=4$  exhaustively. We prohibit repeated blocks, thus extending Kramer's results naturally. This case is of interest for a number of reasons. One is that Billington's methods do not apply here. Perhaps

<sup>1</sup> Research Supported by NSERC Canada under grant A 5047.

<sup>2</sup> Research Supported by NSERC Canada under grant A 7268.

1980 *Mathematics Subject Classification*. Primary 05B05; Secondary 05B30.

*Key words and phrases*. Block design, triple system, decomposable.

more important is the contrast with the cases for smaller  $\lambda$ . For  $\lambda=2$ , every  $B[3, 2; v]$  with  $v \equiv 0, 4 \pmod{6}$  is indecomposable, since  $B[3, 1; v]$  designs do not exist for these orders. Similarly, for  $\lambda=3$ , every  $B[3, 3; v]$  with  $v \equiv 5 \pmod{6}$  is indecomposable. For  $\lambda=4$ , no such trivial indecomposable systems exist. This makes  $\lambda=4$  substantially more complicated. In fact, no examples of indecomposable  $B[3, 4; v]$  designs were known until quite recently. Moreover, it is known that no indecomposable  $NB[3, 4; 7]$  exists [8], and no indecomposable  $NB[3, 4; 9]$  exists [6]. In the remainder of this paper, we establish the following

**MAIN THEOREM.** *An indecomposable  $NB[3, 4; v]$  design exists if and only if  $v \equiv 0, 1 \pmod{3}$  and  $v \geq 10$ .*

## 2. Recursive constructions

Kramer [7] observed that if one could produce a single indecomposable  $NB[3, \lambda; v]$  with  $v$  odd, one could easily produce an infinite family. He based his observation on the following

**LEMMA 2.1** [7, Thm 2.2]: *If there exists an indecomposable  $NB[3, \lambda; v]$  with  $v$  odd there exists an indecomposable  $NB[3, \lambda; 2v+1]$ .* ■

Kramer's proof involves embedding the system of order  $v$  in a system of order  $2v+1$ ; a decomposition of the entire design would necessitate a decomposition of the (indecomposable) subdesign. The only complexity in the embedding is the avoidance of repeated blocks.

In fact, if repeated blocks are allowed, Stern's theorem [9] guarantees that a  $B[3, \lambda; v]$  design can be embedded in a  $B[3, \lambda; w]$  design for every admissible  $w \equiv 2v+1$ . Hence, Stern's theorem somewhat trivializes the existence problem for indecomposable triple systems with repeated blocks. However, no result analogous to Stern's is known when repeated blocks are forbidden.

Fortunately, we do not require such a general result. For our purposes, it suffices to embed a system of order  $v$  into selected larger orders. We do this in a succession of three theorems.

**THEOREM 2.2.** *An  $NB[3, \lambda; v]$  can be embedded in an  $NB[3, \lambda; 2v+1]$ .*

**PROOF.** Lemma 2.1 handles the case when  $v$  is odd. When  $v$  is even, the  $NB[3, \lambda; 2v+1]$  has elements  $\{x_1, \dots, x_v\} \cup \{y_1, \dots, y_{v+1}\}$ . On the  $\{x_i\}$ , place a copy of the  $NB[3, \lambda, v]$ . Since  $v+1$  is odd, the complete graph on the  $v+1$  vertices  $\{y_1, \dots, y_{v+1}\}$  has even degree, and hence has a 2-factorization into 2-factors  $F_1, \dots, F_{v/2}$ . Each 2-factor  $F_i$  contains  $v+1$  edges  $f_{ij}$ ,  $0 \leq j \leq v$ . To complete the design, whenever  $f_{ij} = \{y_a, y_b\}$ , add the  $\lambda/2$  triples  $\{y_a, y_b, x_{i+s}\}$ , and the  $\lambda/2$  triples  $\{y_a, y_b, x_{v/2+i+s}\}$ ,  $0 \leq s < \lambda/2$ , subscripts modulo  $v$ . This is an  $NB[3, \lambda; 2v+1]$  design.

**THEOREM 2.3.** *An  $NB[3, \lambda; v]$  design can be embedded in an  $NB[3, \lambda; 2v+4]$  design provided  $\lambda$  is even.*

**PROOF.** Suppose  $v$  is odd, and write  $\lambda=2t$ . The  $NB[3, \lambda; 2v+4]$  has elements  $\{x_1, \dots, x_v\} \cup \{y_1, \dots, y_{v+4}\}$ . On the  $\{x_i\}$ , place a copy of the  $NB[3, 4; v]$ . There

are two remaining types of triples, those involving a member of  $\{x_i\}$ , and those not. There are  $t(v+4)$  of the latter type, and they are produced cyclically as follows (for definitions on cyclic systems, see [5]). We choose any  $t$  distinct difference triples on the  $\{y_i\}$ , and include the triples obtained by expanding the corresponding starter blocks modulo  $v+4$ . Each starter block covers one occurrence of three differences. Ultimately, each difference must be covered  $t$  times. In order to cover the remainder, note that each difference induces a 2-factor on the  $\{y_i\}$ . With each  $\{x_i\}$ , we therefore associate  $t$  distinct differences  $\{d_{ij}, 0 \leq j < t\}$ , in such a way that each difference appears  $t$  times overall. This can be done in many ways. For example, if the list of differences yet to be used is sorted, we simply set  $d_{ij}$  to be the  $(jv+i)$ 'th element of the list. We then include the triples  $\{x_i, y_j, y_k\}$ ,  $0 \leq j < v$ ,  $k \equiv j + d_{is} \pmod{v+4}$  for some  $s$ . The result is an  $NB[3, \lambda; 2v+4]$ .

When  $v$  is even, the proof is very similar. The only change is that the difference  $(v+4)/2$  induces a 1-factor, not a 2-factor. To avoid this difficulty, we simply ensure that all  $t$  occurrences of this difference are used in the  $t$  starter blocks chosen. ■

**THEOREM 2.4.** *An  $NB[3, \lambda; v]$  can be embedded in an  $NB[3, \lambda; 2v+7]$ .*

**PROOF.** The proof is similar to that of Theorem 2.3. ■

Theorems 2.2 and 2.3 together have the following

**COROLLARY 2.5.** *If there are indecomposable  $NB[3, 4; v]$  for  $v=10, 12, 13, 15, 16, 18, 19$ , and  $22$ , the Main Theorem follows.*

**PROOF.** Consider an arbitrary  $v \geq 10$ ,  $v \equiv 0, 1 \pmod{3}$ , other than one on the list given. If  $v \equiv 0, 4 \pmod{3}$ ,  $(v-4)/2 \equiv 0, 1 \pmod{3}$  and since there is an indecomposable  $NB[3, 4; (v-4)/2]$ , Theorem 2.3 gives an indecomposable  $NB[3, 4; v]$ . If  $v \equiv 1, 3 \pmod{6}$ ,  $(v-1)/2 \equiv 0, 1 \pmod{3}$ , and there is an indecomposable  $NB[3, 4; (v-1)/2]$ ; Theorem 2.2 gives an indecomposable  $NB[3, 4; v]$ . ■

### 3. Small orders

For the eight required orders, we are unable to simply embed a smaller indecomposable design. However, a similar technique does settle three of the cases. Observe that when  $v \equiv 0 \pmod{2}$ , the design can have no decomposition into two systems, one with  $\lambda=1$  and the other with  $\lambda=3$ . Thus any decomposition produces two systems with  $\lambda=2$ . Next, note that the unique  $NB[3, 4; 7]$ , although decomposable, can only be decomposed into  $\lambda=3+1$  (i.e., into an  $NB[3, 3; 7]$  and an  $NB[3, 1; 7]$ ). Thus, any  $B[3, 4; v]$ ,  $v$  even, with the  $NB[3, 4; 7]$  as a sub-design, is indecomposable. With this, Theorem 2.3 immediately produces an indecomposable  $NB[3, 4; 18]$ , and similar constructions settle  $v=16$  and  $v=22$ .

For odd orders, the  $NB[3, 4; 7]$  can still be used to prohibit a  $\lambda=2+2$  decomposition, but another "gadget" is needed to prohibit a  $\lambda=3+1$  decomposition. The unique  $NB[3, 4; 6]$  can be used for this purpose. Embedding the  $NB[3, 4; 7]$  and the  $NB[3, 4; 6]$  simultaneously ensures indecomposability. At first, this appears to be of little consequence, since it does not apply to sufficiently small orders. However, note that if four blocks containing a fixed pair are deleted from the  $NB[3, 4; 7]$ , the resulting partial  $NB[3, 4; 7]$  still has no  $\lambda=2+2$  decomposition. Embedding



this partial  $NB[3, 4; 7]$  along with the  $NB[3, 4; 6]$ , and completing by computer, produced the following indecomposable  $NB[3, 4; 19]$ :

*ABD ABE ABF ABG ACD ACE ACF ACG ADF ADG AEF AEG BDE BDG BEF  
BFG CDE CDF CEG CFG DEF DEG DFG EFG BCH BCI BCN BCO BHI BHN  
BHO BIN BIO BNO CHI CHN CHO CIN CIO CNO HIN HIO HNO INO BDS  
BES BFS BGS BJM BJP BJQ BJR BKM BKP BKQ BKR BLM BLP BLQ BLR  
BMR BPQ CDS CES CFS CGS CJM CJP CJQ CJR CKM CKP CKQ CKR CLM  
CLP CLQ CLR CMR CPQ DHP DHQ DHR DHS DIP DIQ DIR DIS DJL  
DJO DJQ DJR DKL DKM DKO DKP DLM DLO DMN DMO DNP DNQ DNR  
EHP EHQ EHR EHS EIP EIQ EIR EIS EIJ EJM EJO EJP EKO EKP EKQ  
EKR ELM ELN ELO EMN EMO ENQ ENR FHP FHQ FHR FHS FIL FIP  
FIQ FIR FJK FJN FJO FJP FKN FKO FKR FLN FLO FLQ FMO FMP  
FMQ FNS FNR GHJ GHK GHL GHP KIJ GIK GIL GIP GJN GJO GKN GKS  
GLO GLR GMN GMP GMQ GMS GNQ GOQ GOR GPR GQR HJL HMQ HMR  
HMS IJS IKL IMQ IMR IMS JKN JNS JQS JRS KNS KOS KQS LNP LNS  
LPS LOS LRS MNP MOP OPS OQR AHJ AHK AHL AHM AIJ AIK AIL AIM  
AJL AJM AKL AKM HJK HKL IJK ANP ANQ ANR ANS AOP AOQ AOR AOS  
APR APS AQR AQS OPQ ORS PQR PRS*

Order 15 was settled by a similar approach. The same partial  $NB[3, 4; 7]$  was employed to prevent  $\lambda=2+2$  decomposition. To prevent  $\lambda=3+1$  decomposition, observe that the pairs appearing with a specified element form a  $(v-1)$ -vertex 4-regular graph. Moreover, a  $\lambda=3+1$  decomposition induces a 3-factor and a 1-factor in each of these graphs. We therefore included all triples containing one vertex in such a way that its corresponding graph had no 1-factor. The following indecomposable  $NB[3, 4; 15]$  resulted:

*ACD ACE ACF ACG ADF ADG AEF AEG BCD BCE BCF BCG BDE BDG BEF  
BFG CDE CDF CEG CFG DEF DEG DFG EFG ABL ABM ABN ABO ALM  
ALN ALO AMN AMO ANO ADK AFH AEJ AGI AHI AIJ AHJ AHK AIK AJK  
BLO BLM BLN BMN BMO BNO BDI BEH BFJ BGK BHI BHJ BHK BIJ BIK  
BJK CHL CHM CHN CHO CIL CIM CIN CIO CIL CJL CJM CJN CJO CKL CKM  
CKN CKO DHI DHJ DJK DLM DMN DNO DLO DHM DIN DJO DKL DHN  
DIO DJL DKM GHJ GHK GIJ GLM GMN GNO GLO GHL GIM GJN GKO  
GHM GIN GJO GKL FIK FLN FMO FHL FIM FJN FKO FHN FIO FJL FKM  
FHO FIL FJM FKN EIK ELN EMO EHM EHN EHO EIL EIN EIO EIJ EJM  
EJO EKL EKM EKN HIL IJM JKN HKO*

The remaining orders were all settled using sophisticated backtracking; for orders 12 and 13, large portions of the design were first constructed by hand, and the systems were completed by computer. For order 10, much computation was required; over six hundred systems were examined before an indecomposable one was found. The three final orders are given here:

#### Order 10

*ABG ABH ABI ABJ ACG ACH ACI ACJ ADF ADH ADI ADJ AEF AEG AEI  
AEJ AFG AFH BCG BCH BCI BCJ BDF BDG BDI BDJ BEF BEH BEI BEJ BFG  
BFH CDE CDF CDG CDH CEF CEG CEH CFI CFJ CIJ DEF DEI DEJ DGH  
DGJ DHI EGH EGI EHJ FGI FGJ FHI FHJ FIJ GHI GHJ GIJ HIJ*



*Order 12*

ABC ABD ABE ABG ACD ACE ACF AFG ADE ADF AEF AGH AGI AHJ AHK  
 AHL AIJ AIK AIL AJK AJL AKL BCD BCE BCI BDE BDI BEI BGI BFG BGH  
 BFJ BFK BFL BHJ BHK BHL BJK BJL BKL EFI IKL CDK CDL CEK CEL CFJ  
 CFK CFL CGH CGI CGJ CGL CHI CHJ CHK CIJ DEK DEL DFG DFH DFJ  
 DGJ DGK DGL DHI DHJ DHL DIJ DIK EFG EFH EGH EGJ EGK EHI EHL  
 EIJ EJK EKL FHI FHK FIK FIL FJL GIL GJK GKL

*Order 13*

ABC ADE BDE CDE ABF ABG ABH AFG AFH AGH ADF ADG ADH BDF  
 BDG BDH BFG BFH BGH BCI BCJ BCK CDI CDJ CFH CGH ACJ ACK ACL  
 AEJ AEK AEL AIJ AIK AIL AIM AJM AKM ALM BEJ BEK BEL BIK BIL BIM  
 BJL BJM BKM BLM CDM CEK CEL CEM CFK CFL CFM CGI CGJ CGL  
 CHI CHM DEM DFL DFM DGL DGM DHJ DHK DIJ DIK DIL DJK DKL EFH  
 EFI EFJ EFM EGH EGI EGJ EGL EHI EHK EIM FGI FGK FIJ FIK FJK FJL  
 FLM GIM GJK GJM GKL GKM HIJ HIL HJL HJM HKL HKM HLM JKL

This completes the proof of the Main Theorem.

#### 4. Concluding remarks

Our Main Theorem is strong evidence that indecomposability is not merely the result of numerical conditions giving rise to trivial examples. Because of this, it seems reasonable to expect that indecomposable triple systems exist for every  $\lambda$ , with finitely many exceptions for each  $\lambda$ . The existence problem is now settled for  $\lambda=2, 3$  and 4. Infinite families for every odd  $\lambda$  follow from [4]. For  $\lambda=6$ , trivial examples exist, whenever  $v \equiv 2 \pmod{6}$ . Beyond this, nothing is known.

Two avenues of further research seem most appropriate. For fixed small  $\lambda$ , it is necessary to identify features which can be exploited to ensure indecomposability, as we did in Section 3. We are currently studying similar techniques for  $\lambda=5$  and  $\lambda=6$ . The second avenue is to consider the problem for general  $\lambda$ . To supplement recursive constructions such as those in Section 2, a next major step is the production of indecomposable triple systems for every even  $\lambda$ .

**ACKNOWLEDGEMENTS.** We would like to thank Marlene Colbourn and Eric Mendelsohn for their helpful comments.

#### REFERENCES

- [1] BILLINGTON, E. J., Construction of some irreducible designs, *Proceedings of the Ninth Australian Conference on Combinatorial Mathematics*, 1981, 182—196.
- [2] BILLINGTON, E. J., Further constructions of irreducible designs, *Proceedings of the Eleventh Manitoba Conference on Numerical Mathematics and Computing*, 1981, 77—89.
- [3] COLBOURN, C. J. and COLBOURN, M. J., The computational complexity of decomposing block designs, *Ann. Discrete Math.* **26** (1985), 345—350.
- [4] COLBOURN, C. J. and COLBOURN, M. J., Decomposition of block designs: computational issues, *Proceedings of the Tenth Australian Conference on Combinatorial Mathematics*, 1982, 141—146.
- [5] COLBOURN, M. J. and MATHON, R. A., On cyclic Steiner 2-designs, *Ann. Discrete Math.* **7** (1980), 215—253.

- [6] HARNAU, W., Die Anzahl paarweise nichtisomorpher, elementarer, wiederholungsfreier 2-(9, 3, 1)-Blockpläne, *Rostock. Math. Kolloq.* **13** (1980), 43—47.
- [7] KRAMER, E. S., Indecomposable triple systems, *Discrete Math.* **8** (1974), 173—180. *MR* **48** # 10863.
- [8] NETTO, E., *Lehrbuch der Combinatorik*, Teubner, Leipzig, 1927.
- [9] STERN, G., Tripelsysteme mit Untersystemen, *Arch. Math. (Basel)* **33** (1979), 204—208. *MR* **81f**: 05048.

(Received July 27, 1982; in revised form July 20, 1983)

DEPARTMENT OF COMPUTATIONAL SCIENCE  
UNIVERSITY OF SASKATCHEWAN  
SASKATOON, SASKATCHEWAN  
S7N 0W0

DEPARTMENT OF MATHEMATICAL SCIENCES  
MCMASTER UNIVERSITY  
HAMILTON, ONTARIO  
L8S 4K1  
CANADA

## SPLINE INTERPOLATION IN TWO VARIABLES

MARGIT LÉNÁRD

The present paper contains a construction of a two-variable spline function and convergence theorems concerning it. The advantage of this construction is its simplicity on the one hand, and that the sequence of spline functions and their derivatives converge uniformly “in the best way” on the other hand. It means, that the speed of convergence is just the same as that of the best approximating polynomials, which cannot be exceeded at all. Another advantage of the construction is that it rests on the given function values only and there is no need for the values of the derivatives. Further we notice that the present method can be generalized for functions with more than two variables, because the construction is based on the notion of the Taylor formula in several variables.

Without loss of generality we may assume that the domain of the function being approximated is a rectangle in the plane and the values of this function are known at the angular points of a rectangle-lattice, at the so-called “lattice-points”.

Other approaches of this problem can be found in [1], [3], [4], [5] and [6].

As an application we remark that the “net method” (see e.g. [2], [5]) of solving second order partial differential equations gives the approximate values only at the lattice points. By our method we can extend the approximating function from the lattice points to the whole domain.

The following notations are used:

- $f^{(i)}$  denotes the  $i$ -th derivative of the one-variable function  $f$ ;
- $\partial_1^p \partial_2^q u_{j,k}$  denotes the partial derivative of the  $p$ -th order in the first variable and of the  $q$ -th order in the second variable of the two-variable function  $u$  at the lattice-point  $(x_j, y_k)$ ;
- $\Delta^{p,q} u_{j,k}$  denotes the partial difference of the  $p$ -th order in the first variable and of the  $q$ -th order in the second variable of the two-variable function  $u$  at the lattice-point  $(x_j, y_k)$ ;
- $C^n(D)$  denotes the class of  $n$  times continuously differentiable functions on  $D$ ;
- $\omega_n(h)$  denotes the common continuity modulus of the  $n$ -th order partial derivatives of the function  $u$ , that is

$$\omega_n(h) = \max_{\substack{p,q=0,1,\dots,n \\ p+q=n}} \sup_{|x-x'|^2 + |y-y'|^2 \leq h^2} |\partial_1^p \partial_2^q u(x, y) - \partial_1^p \partial_2^q u(x', y')|.$$

Let  $D = \{(x, y): a \leq x \leq b, c \leq y \leq d\}$  a rectangle in the plane. Further let

$$D_\Delta = \{(x_j, y_k): a = x_0 < x_1 < \dots < x_{n+1} = b, c = y_0 < y_1 < \dots < y_{m+1} = d\}$$

---

1980 *Mathematics Subject Classification*. Primary 41A15, 65D07; Secondary 41A05, 65D05.  
*Key words and phrases*. Approximation, interpolation, two variable spline functions.

a subdivision of  $D$  with  $h$  and  $l$ , where  $x_{j+1}-x_j=h$ ,  $y_{k+1}-y_k=l$  ( $j=0, 1, \dots, n$ ,  $k=0, 1, \dots, m$ ). Let

$$D_{h,l} = \{(x, y): x_1 \leq x \leq x_n, y_1 \leq y \leq y_m\}.$$

Let the values  $A_{j,k}^{(p,q)}$  ( $p, q=0, 1, 2$ ,  $p+q \leq 2$ ,  $j=1, \dots, n$ ,  $k=1, \dots, m$ ) be given. We are looking for that spline function  $S_A$  in two variables which satisfies the following conditions:

- (a) it is twice continuously partially differentiable at the lattice-points  $(x_j, y_k)$  ( $j=1, \dots, n$ ,  $k=1, \dots, m$ ) and

$$\partial^p_1 \partial^q_2 S_A(x_j, y_k) = p! q! A_{j,k}^{(p,q)},$$

where  $p, q=0, 1, 2$ ,  $p+q \leq 2$ ;

- (b) on each lattice rectangle it is a polynomial in both variables of minimal degree;  
(c) on each lattice rectangle its degree as regarded a two-variable polynomial is minimal.

If the function  $S_A$  is regarded at fixed  $x=x_j$  or  $y=y_k$ , the condition (a) implies that on each lattice rectangle  $S_A$  is a polynomial in both variables of degree at least five.

Let for  $(x, y) \in [x_j, x_{j+1}] \times [y_k, y_{k+1}]$  ( $j=1, \dots, n-1$ ;  $k=1, \dots, m-1$ )

$$(1) \quad S_A(x, y) = S_{j,k}(x, y) = \sum_{\substack{\mu, \nu=0, 1, \dots, 5 \\ \mu+\nu \leq 6}} A_{j,k}^{(\mu, \nu)} (x-x_j)^\mu (y-y_k)^\nu.$$

If there exist constants  $A_{j,k}^{(\mu, \nu)}$  ( $\mu, \nu=0, 1, \dots, 5$ ,  $3 \leq \mu+\nu \leq 6$ ) for which (1) satisfies condition (a), then (1) satisfies also (b) and (c). The function  $S_A$  is twice continuously partially differentiable at the lattice-points  $(x_j, y_k)$  if and only if the functions  $S_{j,k}$  satisfy

$$(2) \quad \partial^p_1 \partial^q_2 S_{j,k}(x_s, y_t) = p! q! A_{s,t}^{(p,q)}$$

whenever  $p, q=0, 1, 2$ ,  $p+q \leq 2$ ,  $s=j, j+1$ ,  $t=k, k+1$  ( $j=1, \dots, n$ ,  $k=1, \dots, m$ ). This gives 18 conditions for the unknown coefficients  $A_{j,k}^{(\mu, \nu)}$  ( $\mu, \nu=0, 1, \dots, 5$ ,  $3 \leq \mu+\nu \leq 6$ ). The matrix of the system of equations (2) is  $A$ , where the order of unknowns:  $A_{j,k}^{(i,0)}$ ,  $A_{j,k}^{(i-1,i)}$ ,  $\dots$ ,  $A_{j,k}^{(0,i)}$ ,  $i=3, 4, 5, 6$  (without  $A_{j,k}^{(6,0)}$ ,  $A_{j,k}^{(0,6)}$  of course).

As  $\text{Rg } A=18$ , two unknowns may be given arbitrary values, the system of equations can be solved and it has infinite number of solutions. Because we are to minimize the degree it seems suitable to choose the unknowns  $A_{j,k}^{(4,2)}$ ,  $A_{j,k}^{(2,4)}$  freely.

Let the values  $u_{j,k}$  ( $j=0, 1, \dots, n+1$ ,  $k=0, 1, \dots, m+1$ ) be given and for  $j=1, \dots, n$ ,  $k=1, \dots, m$  let

$$A_{j,k}^{(0,0)} = u_{j,k},$$

$$A_{j,k}^{(1,0)} = \frac{1}{2h} [\Delta^{1,0} u_{j,k} + \Delta^{1,0} u_{j-1,k}],$$

$$A_{j,k}^{(0,1)} = \frac{1}{2l} [\Delta^{0,1} u_{j,k} + \Delta^{0,1} u_{j,k-1}],$$

(3)

$$A_{j,k}^{(2,0)} = \frac{1}{2h^2} \Delta^{2,0} u_{j-1,k},$$

$$A_{j,k}^{(1,1)} = \frac{1}{2hl} [\Delta^{1,1} u_{j,k} + \Delta^{1,1} u_{j-1,k-1}],$$

$$A_{j,k}^{(0,2)} = \frac{1}{2l^2} \Delta^{0,2} u_{j,k-1}.$$

Let the unknowns  $A_{j,k}^{(4,2)}$ ,  $A_{j,k}^{(2,4)}$  be given the value 0, that is

$$(4) \quad A_{j,k}^{(4,2)} = A_{j,k}^{(2,4)} = 0 \quad (j = 1, \dots, n-1, k = 1, \dots, m-1).$$

With the values (3), (4) the unique solution of the system of equations (2)

$$A_{j,k}^{(3,0)} = -\frac{3}{2h^3} \Delta^{3,0} u_{j-1,k}$$

$$A_{j,k}^{(2,1)} = \frac{1}{2h^2 l} [\Delta^{2,1} u_{j-1,k-1} - \Delta^{2,2} u_{j-1,k-1}]$$

$$A_{j,k}^{(1,2)} = \frac{1}{2hl^2} [\Delta^{1,2} u_{j-1,k-1} - \Delta^{2,2} u_{j-1,k-1}]$$

$$A_{j,k}^{(0,3)} = -\frac{3}{2l^3} \Delta^{0,3} u_{j,k-1}$$

$$A_{j,k}^{(4,0)} = \frac{5}{2h^4} \Delta^{3,0} u_{j-1,k}$$

$$A_{j,k}^{(3,1)} = \frac{1}{2h^3 l} [\Delta^{2,2} u_{j-1,k-1} - 3\Delta^{3,1} u_{j-1,k}]$$

$$A_{j,k}^{(2,2)} = \frac{1}{2h^2 l^2} [\Delta^{2,2} u_{j,k} + 4\Delta^{2,2} u_{j-1,k-1}]$$

$$A_{j,k}^{(1,3)} = \frac{1}{2hl^3} [\Delta^{2,2} u_{j-1,k-1} - 3\Delta^{1,3} u_{j,k-1}]$$

(5)

$$A_{j,k}^{(0,4)} = \frac{5}{2l^4} \Delta^{0,3} u_{j,k-1}$$

$$A_{j,k}^{(5,0)} = -\frac{1}{h^5} \Delta^{3,0} u_{j-1,k}$$

$$A_{j,k}^{(4,1)} = \frac{5}{2h^4 l} \Delta^{3,1} u_{j-1,k}$$

$$A_{j,k}^{(3,2)} = -\frac{1}{2h^3l^2} [A^{2,2}u_{j,k} + 2A^{2,2}u_{j-1,k-1}]$$

$$A_{j,k}^{(2,3)} = -\frac{1}{2h^2l^3} [A^{2,2}u_{j,k} + 2A^{2,2}u_{j-1,k-1}]$$

$$A_{j,k}^{(1,4)} = \frac{5}{2hl^4} A^{1,3}u_{j,k-1}$$

$$A_{j,k}^{(0,5)} = -\frac{1}{l^5} A^{0,3}u_{j,k-1}$$

$$A_{j,k}^{(5,1)} = -\frac{1}{h^5l} A^{3,1}u_{j-1,k}$$

$$A_{j,k}^{(3,3)} = \frac{1}{2h^3l^3} [A^{2,2}u_{j,k} + A^{2,2}u_{j-1,k-1}]$$

$$A_{j,k}^{(1,5)} = -\frac{1}{hl^5} A^{1,3}u_{j,k-1}.$$

THEOREM 1. Let the values  $u_{j,k}$  ( $j=0, \dots, n+1$ ,  $k=0, \dots, m+1$ ) be given. Then the two-dimensional spline function  $S_\Delta$  of the form (1) with the coefficients (3), (4) and (5) is continuous on the rectangle  $D_{h,l}$ .

PROOF. We have to show only that  $S_\Delta$  is continuous along the sides of the lattice rectangles.

Let now  $x=x_{j+1}$  ( $j=0, 1, \dots, n-1$ ) be fixed and for  $y_k \leq y \leq y_{k+1}$  ( $k=0, 1, \dots, m-1$ ) let

$$f(y) = S_{j,k}(x_{j+1}, y)$$

and

$$g(y) = S_{j+1,k}(x_{j+1}, y).$$

It is obvious that  $f$  and  $g$  are polynomials of degree at most 5. On the other hand the spline function  $S_\Delta$  satisfies the equations (2) which means that

$$f^{(i)}(y_k) = g^{(i)}(y_k)$$

and

$$f^{(i)}(y_{k+1}) = g^{(i)}(y_{k+1})$$

whenever  $i=0, 1, 2$ . Hence  $y_k$  and  $y_{k+1}$  are zeros of the polynomial  $f-g$  of degree at most 5, and their multiplicity is three, and so  $f(y)=g(y)$  ( $y_k \leq y \leq y_{k+1}$ ).

We have seen that  $S_\Delta$  is continuous with respect to its variable if  $x=x_{j+1}$  is fixed. Similarly we get the analogous result for fixed  $y$ .

LEMMA 1. Let  $u \in C^2(D)$  and  $u_{j,k}=u(x_j, y_k)$  ( $j=0, 1, \dots, n+1$ ,  $k=0, 1, \dots, m+1$ ). Then the constants defined in (3) approximate the respective partial deriv-



atives of the function  $u$  at the lattice-points as follows: for  $j=1, \dots, n$ ,  $k=1, \dots, m$

$$|A_{j,k}^{(1,0)} - \partial_1 u_{j,k}| \leq \frac{h}{2} \omega_2(h),$$

$$|A_{j,k}^{(0,1)} - \partial_2 u_{j,k}| \leq \frac{l}{2} \omega_2(l),$$

$$\left| A_{j,k}^{(2,0)} - \frac{1}{2} \partial_1^2 u_{j,k} \right| \leq \frac{1}{2} \omega_2(h),$$

$$|A_{j,k}^{(1,1)} - \partial_1 \partial_2 u_{j,k}| \leq \frac{(h+l)^2}{hl} \omega_2(\sqrt{h^2 + l^2}),$$

$$\left| A_{j,k}^{(0,2)} - \frac{1}{2} \partial_2^2 u_{j,k} \right| \leq \frac{1}{2} \omega_2(l).$$

PROOF. By  $u \in C^2(D)$  for the fixed value  $y_k$  ( $k=1, \dots, m$ ) we can use the second order Taylor polynomial of the function  $x \rightarrow u(x, y_k)$  at the point  $x_j$  ( $j=1, \dots, n$ ). Substituting  $x=x_{j+1}$  and  $x=x_{j-1}$  we obtain

$$\begin{aligned} u_{j+1,k} &= u_{j,k} + \partial_1 u_{j,k} h + \frac{1}{2} \partial_1^2 u(\xi_{j+1}, y_k) h^2, \\ u_{j-1,k} &= u_{j,k} - \partial_1 u_{j,k} h + \frac{1}{2} \partial_1^2 u(\xi_{j-1}, y_k) h^2 \end{aligned} \quad (6)$$

where  $x_{j-1} < \xi_{j-1} < x_j < \xi_{j+1} < x_{j+1}$ , which implies

$$\begin{aligned} |A_{j,k}^{(1,0)} - \partial_1 u_{j,k}| &= \left| \frac{1}{2h} (u_{j+1,k} - u_{j-1,k}) - \partial_1 u_{j,k} \right| = \\ &= \frac{h}{4} |\partial_1^2 u(\xi_{j+1}, y_k) - \partial_1^2 u(\xi_{j-1}, y_k)| \leq \\ &\leq \frac{h}{4} \omega_2(2h) \leq \frac{h}{2} \omega_2(h). \end{aligned}$$

We get similar result changing the role of  $x$  and  $y$ :

$$|A_{j,k}^{(0,1)} - \partial_2 u_{j,k}| \leq \frac{l}{2} \omega_2(l).$$

Adding the equations (6)

$$\begin{aligned} \left| A_{j,k}^{(2,0)} - \frac{1}{2} \partial_1^2 u_{j,k} \right| &= \left| \frac{1}{2h^2} (u_{j+1,k} - 2u_{j,k} + u_{j-1,k}) - \frac{1}{2} \partial_1^2 u_{j,k} \right| \leq \\ &\leq \frac{1}{4} |\partial_1^2 u(\xi_{j+1}, y_k) - \partial_1^2 u_{j,k}| + \frac{1}{4} |\partial_1^2 u(\xi_{j-1}, y_k) - \partial_1^2 u_{j,k}| \leq \frac{1}{2} \omega_2(h). \end{aligned}$$

We get similarly that

$$\left| A_{j,k}^{(0,2)} - \frac{1}{2} \partial_2^2 u_{j,k} \right| \leq \frac{1}{2} \omega_2(l).$$

Using the second order Taylor formula in two variables for the function  $u \in C^2(D)$  and substituting  $x = x_{j+1}$ ,  $y = y_{k+1}$  and  $x = x_{j-1}$ ,  $y = y_{k-1}$ , respectively, we obtain

$$\begin{aligned} u_{j+1,k+1} &= u_{j,k} + h \partial_1 u_{j,k} + l \partial_2 u_{j,k} + \\ &+ \frac{h^2}{2} \partial_1^2 u(\xi_j, \eta_k) + hl \partial_1 \partial_2 u(\xi_j, \eta_k) + \frac{l^2}{2} \partial_2^2 u(\xi_j, \eta_k), \\ (7) \quad u_{j-1,k-1} &= u_{j,k} - h \partial_1 u_{j,k} - l \partial_2 u_{j,k} + \\ &+ \frac{h^2}{2} \partial_1^2 u(\xi_{j-1}, \eta_{k-1}) + hl \partial_1 \partial_2 u(\xi_{j-1}, \eta_{k-1}) + \frac{l^2}{2} \partial_2^2 u(\xi_{j-1}, \eta_{k-1}), \end{aligned}$$

where  $x_j < \xi_j < x_{j+1}$ ,  $y_k < \eta_k < y_{k+1}$ . These equations and the definition of  $A_{j,k}^{(1,1)}$  yield

$$\begin{aligned} A_{j,k}^{(1,1)} &= \frac{u_{j+1,k+1} - 2u_{j,k} + u_{j-1,k-1}}{2hl} - \frac{h}{l} A_{j,k}^{(2,0)} - \frac{l}{h} A_{j,k}^{(0,2)} = \\ &= \frac{h}{4l} [(\partial_1^2 u(\xi_j, \eta_k) - \partial_1^2 u_{j,k}) + (\partial_1^2 u(\xi_{j-1}, \eta_{k-1}) - \partial_1^2 u_{j,k})] + \\ &+ \frac{1}{2} (\partial_1 \partial_2 u(\xi_j, \eta_k) - \partial_1 \partial_2 u_{j,k}) + \frac{1}{2} (\partial_1 \partial_2 u(\xi_{j-1}, \eta_{k-1}) - \partial_1 \partial_2 u_{j,k}) + \\ &+ \frac{l}{4h} [(\partial_2^2 u(\xi_j, \eta_k) - \partial_2^2 u_{j,k}) + (\partial_2^2 u(\xi_{j-1}, \eta_{k-1}) - \partial_2^2 u_{j,k})] + \\ &+ \frac{h}{l} \left( \frac{1}{2} \partial_1^2 u_{j,k} - A_{j,k}^{(2,0)} \right) + \partial_1 \partial_2 u_{j,k} + \frac{l}{h} \left( \frac{1}{2} \partial_2^2 u_{j,k} - A_{j,k}^{(0,2)} \right). \end{aligned}$$

It follows

$$|A_{j,k}^{(1,1)} - \partial_1 \partial_2 u_{j,k}| \leq \frac{(h+l)^2}{hl} \omega_2(\sqrt{h^2 + l^2}),$$

and so the lemma is proved.

**LEMMA 2.** Let  $u \in C^{p+q-1}(D)$ , where  $p, q$  are nonnegative integers,  $p+q \geq 1$ . Then

$$|\Delta^{p,q} u_{j,k}| \leq h^p l^{q-1} \omega_{p+q-1}(l),$$

if  $q \geq 1$ , and

$$|\Delta^{p,q} u_{j,k}| \leq h^{p-1} l^q \omega_{p+q-1}(h),$$

if  $p \geq 1$ .

PROOF. First we assume, that  $p+q>1$  and  $q \geq 1$ . In this case using the Lagrange theorem  $(p+q-1)$ -times we get

$$\begin{aligned}\Delta^{p,q} u_{j,k} &= \Delta^{p-1,q}(u_{j+1,k} - u_{j,k}) = h \Delta^{p-1,q} \partial_1 u(\xi_j^{(1)}, y_k) = \\ &= h^2 \Delta^{p-2,q} \partial_1^2 u(\xi_j^{(2)}, y_k) = \dots = h^p \Delta^{0,q} \partial_1^p u(\xi_j^{(p)}, y_k) = \\ &= h^p l \Delta^{0,q-1} \partial_1^p \partial_2 u(\xi_j^{(p)}, \eta_k^{(1)}) = \dots = \\ &= h^p l^{q-1} [\partial_1^p \partial_2^{q-1} u(\xi_j^{(p)}, \eta_k^{(q-1)} + l) - \partial_1^p \partial_2^{q-1} u(\xi_j^{(p)}, \eta_k^{(q-1)})],\end{aligned}$$

where

$$\xi_j^{(s-1)} < \xi_j^{(s)} < \xi_j^{(s-1)} + h \quad (s = 1, \dots, p),$$

$$\eta_k^{(t-1)} < \eta_k^{(t)} < \eta_k^{(t-1)} + l \quad (t = 1, \dots, q-1).$$

So obviously

$$|\Delta^{p,q} u_{j,k}| \leq h^p l^{q-1} \omega_{p+q-1}(l).$$

Similarly we get in the case  $p \geq 1$

$$|\Delta^{p,q} u_{j,k}| \leq h^{p-1} l^q \omega_{p+q-1}(h).$$

If  $p+q=1$ , then  $p=1, q=0$  or  $p=0, q=1$ . In the case  $p=1, q=0$

$$|\Delta^{1,0} u_{j,k}| = |u_{j+1,k} - u_{j,k}| \leq \omega(h),$$

and if  $p=0, q=1$

$$|\Delta^{0,1} u_{j,k}| = |u_{j,k+1} - u_{j,k}| \leq \omega(l).$$

So the lemma is proved.

THEOREM 2. Let  $u \in C^2(D)$ ,  $D_\Delta$  be a subdivision of  $D$  with  $h$  and  $l$ ,  $0 < \alpha < h/l < \beta$  and denote  $\|h\| = \sqrt{h^2 + l^2}$ . Then the twodimensional spline function  $S_\Delta$  of the form (1) with the coefficients (3), (4) and (5) approximates the function  $u$  and its derivatives on  $D_{h,l}$  as follows;

$$|\partial_1^p \partial_2^q u(x, y) - \partial_1^p \partial_2^q S_\Delta(x, y)| \leq c_{p,q} \|h\|^{2-(p+q)} \omega_2(\|h\|)$$

for  $p, q=0, 1, 2$ ,  $p+q \leq 2$ , where the constants  $c_{p,q}$  are independent on  $\|h\|$ .

PROOF. Let  $x_j \leq x \leq x_{j+1}$ ,  $y_k \leq y \leq y_{k+1}$  ( $j=1, \dots, n-1$ ,  $k=1, \dots, m-1$ ) and let define

$$\begin{aligned}u^*(x, y) &= u_{j,k} + \partial_1 u_{j,k}(x - x_j) + \partial_2 u_{j,k}(y - y_k) + \frac{1}{2} \partial_1^2 u_{j,k}(x - x_j)^2 + \\ &+ \partial_1 \partial_2 u_{j,k}(x - x_j)(y - y_k) + \frac{1}{2} \partial_2^2 u_{j,k}(y - y_k)^2.\end{aligned}$$

As  $u \in C^2(D)$ , using the second order Taylor formula at the point  $(x_j, y_k)$  for the function  $u$  we have

$$\begin{aligned} |u(x, y) - u^*(x, y)| &\leq \frac{1}{2} |\partial_1^2 u(\xi_j, \eta_k) - \partial_1^2 u_{j,k}| |x - x_j|^2 + \\ &+ |\partial_1 \partial_2 u(\xi_j, \eta_k) - \partial_1 \partial_2 u_{j,k}| |x - x_j| |y - y_k| + \frac{1}{2} |\partial_2^2 u(\xi_j, \eta_k) - \partial_2^2 u_{j,k}| |y - y_k|^2 \leq \\ &\leq \frac{1}{2} (h+l)^2 \omega_2(\sqrt{h^2 + l^2}) \leq \|h\|^2 \omega_2(\sqrt{h^2 + l^2}), \end{aligned}$$

where  $x_j < \xi_j < x$  and  $y_k < \eta_k < y$ . On the other hand by Lemmas 1 and 2 we obtain

$$\begin{aligned} |u^*(x, y) - S_{j,k}(x, y)| &\leq |\partial_1 u_{j,k} - A_{j,k}^{(1,0)}| |x - x_j| + |\partial_2 u_{j,k} - A_{j,k}^{(0,1)}| |y - y_k| + \\ &+ \left| \frac{1}{2} \partial_1^2 u_{j,k} - A_{j,k}^{(2,0)} \right| |x - x_j|^2 + |\partial_1 \partial_2 u_{j,k} - A_{j,k}^{(1,1)}| |x - x_j| |y - y_k| + \\ &+ \left| \frac{1}{2} \partial_2^2 u_{j,k} - A_{j,k}^{(0,2)} \right| |y - y_k|^2 + \sum_{\substack{\mu, \nu=0,1,\dots,6 \\ 3 \leq \mu + \nu \leq 6}} |A_{j,k}^{(\mu,\nu)}| |x - x_j|^\mu |y - y_k|^\nu \leq \\ &\leq h^2 \omega_2(h) + l^2 \omega_2(l) + (h+l)^2 \omega_2(\sqrt{h^2 + l^2}) + 24h^2 \omega_2(h) + 24l^2 \omega_2(l) \leq \\ &\leq 27 \|h\|^2 \omega_2(\|h\|). \end{aligned}$$

Combining these results

$$\begin{aligned} |u(x, y) - S_A(x, y)| &= |u(x, y) - S_{j,k}(x, y)| \leq \\ &\leq |u(x, y) - u^*(x, y)| + |u^*(x, y) - S_{j,k}(x, y)| \leq \\ &\leq 28 \|h\|^2 \omega_2(\|h\|) \end{aligned}$$

and so  $C_{0,0} = 28$ . As  $u \in C^2(D)$ , we get  $\partial_1 u \in C^1(D)$  and so the first order Taylor formula at the point  $(x_j, y_k)$  gives

$$\partial_1 u(x, y) = \partial_1 u_{j,k} + \partial_1^2 u(\xi_j, \eta_k)(x - x_j) + \partial_1 \partial_2 u(\xi_j, \eta_k)(y - y_k),$$

where  $x_j < \xi_j < x$ ,  $y_k < \eta_k < y$ . Let now

$$(\partial_1 u)^*(x, y) = \partial_1 u_{j,k} + \partial_1^2 u_{j,k}(x - x_j) + \partial_1 \partial_2 u_{j,k}(x - x_j)(y - y_k),$$

then

$$\begin{aligned} |\partial_1 u(x, y) - (\partial_1 u)^*(x, y)| &\leq \\ &\leq |\partial_1^2 u(\xi_j, \eta_k) - \partial_1^2 u_{j,k}| |x - x_j| + |\partial_1 \partial_2 u(\xi_j, \eta_k) - \partial_1 \partial_2 u_{j,k}| |y - y_k| \leq \\ &\leq (h+l) \omega_2(\sqrt{h^2 + l^2}) \leq 2 \|h\| \omega_2(\|h\|). \end{aligned}$$

Further

$$\begin{aligned}
 |(\partial_1 u)^*(x, y) - \partial_1 S_{j,k}(x, y)| &\leq |\partial_1 u_{j,k} - A_{j,k}^{(1,0)}| + |\partial_1^2 u_{j,k} - 2A_{j,k}^{(2,0)}| |x - x_j| + \\
 &+ |\partial_1 \partial_2 u_{j,k} - A_{j,k}^{(1,1)}| |y - y_k| + \sum_{\substack{\mu=1, \dots, 5 \\ \nu=0, \dots, 5 \\ 3 \leq \mu + \nu \leq 6}} \mu |A_{j,k}^{(\mu, \nu)}| |x - x_j|^{\mu-1} |y - y_k|^\nu \leq \\
 &\leq \frac{3}{2} h \omega_2(h) + \frac{(h+l)^2}{h} \omega_2(\sqrt{h^2 + l^2}) + \frac{163}{2} h \omega_2(h) + \frac{53}{2} l \omega_2(l) \leq \\
 &\leq \left(113 + \frac{1}{\alpha}\right) \|h\| \omega_2(h).
 \end{aligned}$$

Combining these estimations

$$\begin{aligned}
 |\partial_1 u(x, y) - \partial_1 S_d(x, y)| &\leq |\partial_1 u(x, y) - (\partial_1 u)^*(x, y)| + |(\partial_1 u)^*(x, y) - S_{j,k}(x, y)| \leq \\
 &\leq \left(115 + \frac{1}{\alpha}\right) \|h\| \omega_2(\|h\|),
 \end{aligned}$$

that is  $C_{1,0} = 115 + \frac{1}{\alpha}$ . We obtain similarly for the partial derivative with respect to the second variable by the symmetry  $C_{0,1} = 115 + \beta$ . Now let us consider the approximation of the second order partial derivatives:

$$\begin{aligned}
 |\partial_1^2 u(x, y) - \partial_1^2 S_d(x, y)| &\leq |\partial_1^2 u(x, y) - \partial_1^2 u_{j,k}| + |\partial_1^2 u_{j,k} - \partial_1^2 S_{j,k}(x, y)| \leq \\
 &\leq \omega_2(\sqrt{h^2 + l^2}) + |\partial_1^2 u_{j,k} - 2A_{j,k}^{(2,0)}| + \sum_{\substack{\mu=2, \dots, 5 \\ \nu=0, \dots, 5 \\ 3 \leq \mu + \nu \leq 5}} \mu(\mu-1) |A_{j,k}^{(\mu, \nu)}| |x - x_j|^{\mu-2} |y - y_k|^\nu \leq \\
 &\leq \omega_2(\sqrt{h^2 + l^2}) + \omega_2(h) + 232\omega_2(h) \leq 234\omega_2(\|h\|),
 \end{aligned}$$

that is  $C_{0,2} = C_{2,0} = 243$  (by symmetry). Finally

$$\begin{aligned}
 |\partial_1 \partial_2 u(x, y) - \partial_1 \partial_2 S_d(x, y)| &\leq |\partial_1 \partial_2 u(x, y) - \partial_1 \partial_2 u_{j,k}| + |\partial_1 \partial_2 u_{j,k} - A_{j,k}^{(1,1)}| + \\
 &+ \sum_{\substack{\mu, \nu=1, \dots, 5 \\ 3 \leq \mu + \nu \leq 6}} \mu \nu |A_{j,k}^{(\mu, \nu)}| |x - x_j|^{\mu-1} |y - y_k|^{\nu-1} \leq \\
 &\leq \omega_2(\sqrt{h^2 + l^2}) + \frac{(h+l)^2}{hl} \omega_2(\sqrt{h^2 + l^2}) + 82\omega_2(h) + 82\omega_2(l) \leq \\
 &\leq \left(167 + \frac{1}{\alpha} + \beta\right) \omega_2(\|h\|),
 \end{aligned}$$

that is  $C_{1,1} = 167 + \frac{1}{\alpha} + \beta$ , and so the theorem is proved.

REMARK. On this theorem we have used the following notation

$$\partial_1^p \partial_2^q S_d(x, y) = \partial_1^p \partial_2^q S_{j,k}(x, y)$$

for  $(x, y) \in [x_j, x_{j+1}] \times [y_k, y_{k+1}]$  ( $j=0, 1, \dots, n$ ,  $k=0, 1, \dots, m$ ).

We are very grateful to Professor János Balázs for many valuable discussions and encouragement in the preparation of this work.

## REFERENCES

- [1] AHLBERG, J. H., NILSON, E. N. and WALSH, J. L., *The theory of splines and their applications*, Academic Press, New-York—London, 1967. *MR* 39 #684.
- [2] BEREZIN, I. S. and ZHIDKOV, N. P., *Computing Methods*, Pergamon Press, Oxford—London, 1965. *MR* 30 #4372.
- [3] GRZANNA, J., Zweidimensionale Splineinterpolation über einem Polargitter, *J. Approximation Theory* 22 (1978), 189—201. *MR* 57 #17083.
- [4] Завьялов, Ю. С., Квасов, Б. И. и Мирошниченко, В. Л., *Методы сплайн-функций*, Наука, Москва, 1980.
- [5] Марчук, Г. И., *Методы вычислительной математики*, Наука, Новосибирск, 1973. *MR* 49 #4195.
- [6] Стечкин, С. Б. и Субботин, Ю. Н., *Сплайны в вычислительной математике*, Наука, Москва, 1976. *MR* 56 #13517.

(Received September 16, 1982)

KOSSUTH LAJOS TUDOMÁNYEGYETEM  
P.O. BOX 12  
H-4010 DEBRECEN  
HUNGARY



## BRANCHES CONTINUES DE VECTEURS PROPRES GÉNÉRALISÉS. APPLICATIONS AUX ÉQUATIONS DE COÏNCIDENCES

G. ISAC

1. Les dernières années plusieurs auteurs ont étudié les valeurs propres généralisées, c'est-à-dire pour des couples d'opérateurs  $(f, g)$ , (le nombre  $\lambda \in \mathbf{R}$  est une valeur propre pour le couple  $(f, g)$  où  $f, g: E \rightarrow F$  s'il existe  $x \neq 0$ ,  $x \in E$  tel que  $f(x) = \lambda g(x)$ ) [2], [3], [6], [15], [18], [23], [19], [22], [24], [31].

La raison est que plusieurs équations non-linéaires de la technique ou de la physique mathématique peuvent être présentées sous la forme,  $f(x) = \lambda g(x)$ , qui est une équation de coïncidence à valeurs propres.

Les équations de coïncidence dans des espaces de Banach ont été étudiées par la technique du degré de coïncidence par J. Mawhin [17], R. E. Gaines et J. Mawhin [7], M. Mininni [19] et par la technique du degré topologique par W. V. Petryshyn [23], P. M. Fitzpatrick et W. V. Petryshyn [4], [5], [6].

Des techniques différentes dans l'étude des équations de coïncidence ont été utilisées par S. Fučík, J. Nečas, J. Souček et V. Souček [31] par M. Schechter, J. Shapiro et M. Snow [29] et par E. T. Dean et P. L. Chambre [3].

Dans cet ouvrage nous n'étudions pas le simple problème d'existence des valeurs propres pour le couple  $(f, g)$  mais un problème plus précis.

Etant donnés deux espaces de Fréchet  $E, F$  deux cônes,  $K_1 \subset E$ ,  $K_2 \subset F$ , les opérateurs  $f, g: K_1 \rightarrow K_2$  et deux voisinages de zéro  $U_1, U_2$  dans l'espace  $E$  tels que  $\bar{U}_2 \subset U_1$  et  $U_1 \cap K_1$  est borné, dans quelles hypothèses, étant donné un nombre réel  $\lambda^* > 0$ , il existe  $x^* \in (U_1 \setminus \bar{U}_2) \cap K_1$  tel que,  $f(x^*) = \lambda^* g(x^*)$ ?

Dans l'étude de ce problème nous n'utilisons ni le degré topologique, ni le degré de coïncidence, mais la notion de branche continue de longueur maximale de vecteurs propres pour des couples d'opérateurs. Cette notion a été définie dans le cas particulier des couples  $(f, I)$  par M. A. Krasnoselskii [14], [13], [12], mais qui n'a pas été suffisamment exploitée en analyse non-linéaire.

Ce concept nous permet de donner pour nos résultats des démonstrations topologiques élémentaires.

Nous présentons à la fin quelques applications et quelques directions de futures applications.

Nous remarquons aussi que le théorème 2 est une généralisation pour des couples d'opérateurs  $(f, g)$  du théorème de compression du cône de Krasnoselskii et aussi ce théorème contient comme un cas particulier le théorème classique des valeurs intermédiaires pour une fonction réelle.

1980 *Mathematics Subject Classification*. Primary 47H10; Secondary 46A40.

*Key words and phrases*. Nonlinear operators on convex cones in Fréchet spaces, generalized eigenvalues, coincidence equations on convex cones.

Cet ouvrage est un développement substantiel de notre ouvrage [10].

2. On présente dans cet ouvrage les principaux résultats sur leur forme la plus générale possible donc, on considère des espaces localement convexes.

Si  $E(\tau)$  est un espace localement convexe on suppose que la topologie  $\tau$  est définie par une famille suffisante de semi-normes  $\{\|\cdot\|_\alpha\}_{\alpha \in A}$  [8].

On considère les cônes convexes localement bornés qui ont été utilisés dans les ouvrages [8]—[10].

Soit  $K \subset E$  un cône convexe, c'est-à-dire,  $K + K \subset K$  et  $(\forall \lambda \in \mathbb{R}_+)(\lambda K \subset K)$ ; on dit que le cône  $K$  est *localement borné* [8] si et seulement si, il existe un voisinage  $U$  de zéro tel que  $U \cap K$  est borné.

La notion d'ensemble séparant caractérise les cônes localement bornés; le cône convexe  $K \subset E$  est localement borné si et seulement si, il existe un ensemble  $A \subset K$  séparant c'est-à-dire,  $A$  est fermé dans  $K$ ,  $0 \notin A$  et la composante connexe de  $K \setminus A$  contenant l'origine est un ensemble borné dans  $E$ .

Un cône convexe bien basé [8], [9] est localement borné et donc un cône convexe localement compact ou faiblement localement compact est localement borné [8].

Les cônes convexes semi-complets au sens de Mokobodzki [8] sous certaines hypothèses sont localement bornés [8], [9].

REMARQUE. Chaque cône convexe dans un espace normé est localement borné et donc les résultats qui sont vraies dans les cônes localement bornés sont vraies dans les cônes des espaces normés.

Dans cet ouvrage on utilise la caractérisation suivante d'un cône localement borné [8].

Le cône convexe  $K \subset E$  est *localement borné* si et seulement si, il existe une semi-norme  $p$  continue et uniformément positive sur  $K$  c'est-à-dire:

$$(\forall \alpha \in \mathcal{A})(\exists \gamma_\alpha > 0)(\forall x \in K)(\gamma_\alpha |x|_\alpha \leq p(x)).$$

3. Soit  $E(\tau)$  un espace localement convexe et  $\mathcal{C}$  un cône convexe abstrait muni d'une structure de Riesz, c'est-à-dire on suppose qu'il existe sur  $\mathcal{C}$  une structure d'espace réticulé compatible avec la structure de cône convexe donnée sur  $\mathcal{C}$ .

Soit  $\mathcal{B}_E$  la famille des ensembles bornés de l'espace  $E$ .

On dit que la fonction  $\alpha: \mathcal{B}_E \rightarrow \mathcal{C}$  est une *mesure de non-compacité* sur  $E$  si les propriétés suivantes sont vérifiées:

$$1^\circ) (\forall \Omega_1, \Omega_2 \in \mathcal{B}_E)(\Omega_1 \subset \Omega_2) \Rightarrow (\alpha(\Omega_1) \leq \alpha(\Omega_2))$$

$$2^\circ) (\forall \Omega \in \mathcal{B}_E)(\alpha(\overline{\text{co}}(\Omega)) = \alpha(\Omega))$$

$$3^\circ) (\forall \Omega_1, \Omega_2 \in \mathcal{B}_E)(\alpha(\Omega_1 \cup \Omega_2) = \sup \{\alpha(\Omega_1), \alpha(\Omega_2)\})$$

$$4^\circ) (\forall \Omega_1, \Omega_2 \in \mathcal{B}_E)(\alpha(\Omega_1 + \Omega_2) \leq \alpha(\Omega_1) + \alpha(\Omega_2))$$

$$5^\circ) (\forall \Omega \in \mathcal{B}_E)(\forall \lambda \in \mathbb{R})(\alpha(\lambda \cdot \Omega) = |\lambda| \alpha(\Omega))$$

$$6^\circ) (\forall \Omega \in \mathcal{B}_E)(\alpha(\Omega) = 0 \Leftrightarrow \Omega \text{ est précompact}).$$

EXEMPLES. 1°) La mesure de non-compacité de Hausdorff [26].

2°) Soit  $E(\tau)$  un espace de Fréchet qui a la topologie  $\tau$  définie par la famille suffisante de semi-normes  $\{\|_n\}_{n \in \mathbb{N}}$ . On prend comme cône  $\mathcal{C}$  l'ensemble de toutes les fonctions  $a: \mathbb{N} \rightarrow [0, +\infty[$  et on considère sur  $\mathcal{C}$  la structure naturelle de cône convexe donnée par :  $(a_1 + a_2)(n) = a_1(n) + a_2(n)$  pour tout  $a_1, a_2 \in \mathcal{C}$  et tout  $n \in \mathbb{N}$  et  $(\lambda \cdot a)(n) = \lambda \cdot a(n)$  quel que soit  $\lambda \in \mathbb{R}_+$ ,  $a \in \mathcal{C}$  et  $n \in \mathbb{N}$ . La relation d'ordre:  $a_1 \leq a_2 \Leftrightarrow (\forall n \in \mathbb{N}) (a_1(n) \leq a_2(n))$  donne sur  $\mathcal{C}$  une structure d'espace de Riesz. La fonction  $\alpha: \mathcal{B}_E \rightarrow \mathcal{C}$  définie par :

$$[\alpha(\Omega)](n) = \inf \{d > 0 \mid \Omega = \bigcup_{i=1}^m \Omega_i, \|_n - \text{diam.}(\Omega_i) \leq d\}$$

vérifie les propriétés 1°)—6°) et elle s'appelle *la mesure de non-compacité de Kuratowski*.

3°) Dans les ouvrages [26], [25], [4] on trouve plusieurs exemples de mesures de non-compacité.

Soit  $E(\tau_1)$ ,  $F(\tau_2)$  deux espaces localement convexes,  $K_1 \subset E$ ,  $K_2 \subset F$  deux cônes convexes et  $\mathcal{C}$  un cône convexe abstrait muni d'une structure de Riesz.

Soit  $\alpha_1: \mathcal{B}_E \rightarrow \mathcal{C}$ ,  $\alpha_2: \mathcal{B}_F \rightarrow \mathcal{C}$  deux mesures de non-compacité précisées.

On dit que l'opérateur  $f: K_1 \rightarrow K_2$  est de type  $(k, \alpha_1, \alpha_2)$ -Lipschitz si les affirmations suivantes sont vérifiées :

- 1<sub>1</sub>)  $(\forall \Omega \in \mathcal{B}_{K_1}) (f(\Omega) \in \mathcal{B}_{K_2})$
- 1<sub>2</sub>)  $f$  est continu
- 1<sub>3</sub>)  $(\exists k \in \mathbb{R}_+) (\forall \Omega \in \mathcal{B}_{K_1}) (\alpha_2(f(\Omega)) \leq k\alpha_1(\Omega))$

où  $\mathcal{B}_{K_1}$  (resp.  $\mathcal{B}_{K_2}$ ) est la famille des ensembles bornés du cône  $K_1$  (resp.  $K_2$ )

REMARQUE. Si  $F$  est un espace localement convexe quasi-complet et  $f, (k, \alpha_1, \alpha_2)$ -Lipschitz où  $k=0$  alors  $f$  est complètement continu.

Exemples d'opérateurs  $(k, \alpha_1, \alpha_2)$ -Lipschitz se trouvent dans les ouvrages [4], [25], [26].

4. Soit  $E(\tau_1)$ ,  $F(\tau_2)$  deux espaces localement convexes,  $K_1 \subset E$ ,  $K_2 \subset F$  deux cônes convexes et  $f, g: K_1 \rightarrow K_2$  deux opérateurs.

On considère les notions de vecteurs propres et de valeurs propres au sens généralisé comme dans l'ouvrage [31] pour le couple  $(f, g)$ .

On pose :

$$\mathcal{S}(f, g) = \{x \in K_1 \setminus \{0\} \mid \exists \lambda \geq 0, f(x) = \lambda g(x)\}.$$

DÉFINITION 1. Si  $K_1$  est un cône convexe localement borné on dit que  $\mathcal{S}(f, g)$  est une *branche continue de longueur maximale* de vecteurs propres sur l'ensemble  $A \subset K_1$  si, pour tout voisinage ouvert  $U$  de zéro tel que  $U \cap K_1$  est borné et  $\partial U \cap A \neq \emptyset$  on a,  $\partial U \cap \mathcal{S}(f, g) \cap A \neq \emptyset$ .

Pour les résultats qu'on présente dans cet ouvrage il est important de savoir quelques critères d'existence des branches continues de vecteurs propres.

Soit  $E$  un espace de Banach et  $E^*$  son dual topologique.

On note par  $<, >$  l'application bilinéaire de dualité,  $\rightharpoonup$  la convergence faible et par  $\rightarrow$  la convergence forte.

Soit  $K_1 \subset E$ ,  $K_2 \subset E^*$  deux cônes convexes fermés.

On dit que l'opérateur  $f: K_1 \rightarrow K_2$  vérifie la condition (S) si, pour toute suite  $\{x_n\}_{n \in \mathbb{N}} \subset K_1$  qui a la propriété que  $\{x_n\}_{n \in \mathbb{N}} \rightarrow x_0$  (dans le cône  $K_1$ ) et  $\langle f(x_n) - f(x_0), x_n - x_0 \rangle \rightarrow 0$  il résulte que  $\{x_n\}_{n \in \mathbb{N}} \rightarrow x_0$ .

On dit que l'opérateur  $f: K_1 \rightarrow K_2$  vérifie la condition (S<sub>0</sub>) si, pour toute suite  $\{x_n\}_{n \in \mathbb{N}} \subset K_1$  telle que  $\{x_n\}_{n \in \mathbb{N}} \rightarrow x_0$  (dans le cône  $K_1$ ),  $\{f(x_n)\}_{n \in \mathbb{N}} \rightarrow y \in K_2$  et  $\langle f(x_n), x_n \rangle \rightarrow \langle y, x_0 \rangle$  il résulte que  $\{x_n\}_{n \in \mathbb{N}} \rightarrow x_0$ .

REMARQUE. Si l'opérateur  $f$  vérifie la condition (S) alors il vérifie la condition (S<sub>0</sub>).

PROPOSITION 1. Soit  $E$  un espace de Banach séparable et réflexif,  $K_1 \subset E$ ,  $K_2 \subset E^*$  deux cônes convexes fermés,  $f: K_1 \rightarrow K_2$  un opérateur continu, borné, homogène de degré  $r$  qui vérifie la condition (S<sub>0</sub>) et  $g: K_1 \rightarrow K_2$  un opérateur complètement continu homogène de degré  $s$ .

On suppose qu'il est fixée une suite  $\{E_n\}_{n \in \mathbb{N}}$  croissante de sous-espaces  $E_n \subset E$  tels que,  $\dim E_n = n$  et  $\bigcup E_n = E$ . Soit  $B \subset K_1$  un sous-ensemble borné et fermé tel que  $0 \notin B$ .

Soit  $I_n: E_n \rightarrow E$  l'injection canonique et  $P_n: E^* \rightarrow E_n^*$  la projection duale.

Soit  $K_{1n} = K_1 \cap E_n$ , (on suppose  $K_{1n} \neq \{0\}$ ),  $K_{2n} = P_n(K_2)$ , (on suppose  $K_{2n} \neq \{0\}$ ) pour tout  $n \in \mathbb{N}$ .

Si, pour chaque  $n \in \mathbb{N}$  il existe  $\lambda_n > 0$  et  $x_n \in B \cap K_{1n}$  tel que,  $P_n f(x_n) = \lambda_n P_n g(x_n)$  et la suite  $\{\lambda_n\}_{n \in \mathbb{N}}$  est bornée, alors il existe une branche continue de longueur maximale de vecteurs propres pour le couple  $(f, g)$  sur le cône  $K_1$ .

DÉMONSTRATION. D'abord on montre qu'il existe  $x_0 \in B$  et  $\lambda_0 > 0$  tels que  $f(x_0) = \lambda_0 g(x_0)$ .

Puisque  $B$  est borné et  $E$  réflexif il résulte que la suite  $\{x_n\}_{n \in \mathbb{N}}$  contient une sous-suite faiblement convergente. On peut supposer (éventuellement on considère une sous-suite) que cette sous-suite est même  $\{x_n\}_{n \in \mathbb{N}}$  et soit donc,  $\{x_n\}_{n \in \mathbb{N}} \rightarrow x_0 \in K_1$ .

On prouve maintenant que la suite  $\{x_n\}_{n \in \mathbb{N}}$  contient une sous-suite convergente en norme vers  $x_0$  et que  $x_0$  est un vecteur propre pour  $(f, g)$ .

Puisque  $\{\lambda_n\}_{n \in \mathbb{N}}$  est bornée (éventuellement on considère une sous-suite) on peut supposer qu'il existe  $0 < \lambda_0 = \lim_{n \rightarrow \infty} \lambda_n$ .

L'opérateur  $g$  étant complètement continu on peut supposer (éventuellement on considère une sous-suite) que  $\{g(x_n)\}_{n \in \mathbb{N}} \rightarrow w \in K_2$ .

Soit  $m \in \mathbb{N}$  fixé et  $v \in E_m$ ; si  $n \geq m$  ( $E_m \subset E_n$ ) alors,  $\langle f(x_n), v \rangle = \langle f(x_n), I_n(v) \rangle = \langle P_n(f(x_n)), v \rangle = \lambda_n \langle P_n(g(x_n)), v \rangle = \lambda_n \langle g(x_n), v \rangle = \langle \lambda_n g(x_n), v \rangle$  d'où il résulte que  $\langle f(x_n), v \rangle \rightarrow \langle \lambda_0 w, v \rangle$ .

La dernière relation est vraie pour tout  $v \in \bigcup_{n \in \mathbb{N}} E_n$ .

Puisque  $\bigcup E_n = E$  et l'ensemble  $\{f(x_n)\}_{n \in \mathbb{N}}$  est borné on peut prouver que  $\{f(x_n)\}_{n \in \mathbb{N}} \rightarrow \lambda_0 w$ .

On a aussi que,  $\langle f(x_n), x_n \rangle = \langle \lambda_n g(x_n), x_n \rangle = \lambda_n \langle g(x_n), x_n \rangle$  et on peut prouver que  $\langle f(x_n), x_n \rangle \rightarrow \langle \lambda_0 w, x_0 \rangle$ .

Utilisant maintenant la condition (S<sub>0</sub>) pour  $f$  on obtient que  $\{x_n\}_{n \in \mathbb{N}} \rightarrow x_0$  et comme  $f$  et  $g$  sont continus on a,  $f(x_0) = \lim_{n \rightarrow \infty} f(x_n) = \lambda_0 w = \lambda_0 g(x_0)$ .

Puisque l'ensemble  $B$  est fermé et  $0 \notin B$  il résulte que  $x_0 \neq 0$ .

Soit maintenant  $U$  un voisinage de zéro de l'espace  $E$  tel que  $U \cap K_1$  est borné.

Il existe  $\varrho_0 > 0$  tel que  $\varrho_0 x_0 \in \partial U \cap K_1$  et alors on a,  $\lambda_0 g(\varrho_0 x_0) = \lambda_0 \varrho_0^s g(x_0) = \varrho_0^s f(x_0) = \varrho_0^s \left[ f \left( \frac{1}{\varrho_0} (\varrho_0 x_0) \right) \right] = \varrho_0^{s-r} f(\varrho_0 x_0)$ , d'où,  $f(\varrho_0 x_0) = \lambda_0 \varrho_0^{r-s} g(\varrho_0 x_0)$ .

Donc le couple  $(f, g)$  a une branche continue de longueur maximale de vecteurs propres sur le cône  $K_1$ .  $\square$

On considère maintenant le cas particulier des couples de la forme  $(f, I)$  où  $I$  est l'opérateur identique.

Si on analyse la démonstration du théorème 1 de l'ouvrage [9] on observe que le théorème est vraie même pour  $k=0$ .

Donc on a le théorème suivant.

**THÉORÈME 1.** Soit  $E(\tau)$  un espace de Fréchet,  $K \subset E$  un cône convexe fermé, normal, saillant et  $f: K \rightarrow K$  un opérateur  $(k, \alpha, \alpha)$ -Lipschitz ( $k \geq 0$ ).

On suppose qu'il existe un ensemble  $A \subset K$  séparant.

Alors:

1°) il existe une semi-norme  $p$  continue, monotone croissante et uniformément positive sur  $K$ .

2°) si:  $\delta_A > k d_A$  où:  $\delta_A = \inf \{p(f(x)) | x \in A\}$  et  $d_A = \sup \{p(x) | x \in A\}$ , alors il existe  $\lambda > 0$  et  $z \in A$  tels que  $\lambda f(z) = z$ .  $\square$

Utilisant ce théorème et l'ouvrage [9] on obtient les résultats suivants.

**PROPOSITION 2.** Soit  $E(\tau)$  un espace de Fréchet,  $K \subset E$  un cône convexe fermé, saillant, normal, localement borné et  $p$  une semi-norme continue monotone croissante et uniformément positive sur  $K$ .

Si  $f: K \rightarrow K$  est un opérateur  $(k, \alpha, \alpha)$ -Lipschitz où  $k > 0$  et  $\sup \{\lambda | (x, \lambda) \in \mathcal{C}\} < \frac{1}{k}$  où  $\mathcal{C}$  est la composante connexe de l'ensemble  $S = \{(x, \lambda) \in K \times \mathbb{R}_+ | x = \lambda f(x)\}$  qui contient  $(0, 0)$ , alors le couple  $(f, I)$  a une branche continue de longueur maximale de vecteurs propres sur le cône  $K$ .

**PROPOSITION 3.** Soit  $E(\tau)$  un espace de Fréchet,  $K \subset E$  un cône convexe, fermé, saillant, localement borné et  $f: K \rightarrow K$  un opérateur complètement continu.

Si  $S = \{(x, \lambda) \in K \times \mathbb{R}_+ | x = \lambda f(x)\}$  et  $\mathcal{C}$  est la composante connexe de  $S$  qui contient  $(0, 0)$  et si,  $\sup \{\lambda | (x, \lambda) \in \mathcal{C}\} \neq +\infty$  alors le couple  $(f, I)$  a une branche continue de longueur maximale de vecteurs propres sur le cône  $K$ .  $\square$

La résultat suivant est une généralisation du critère du minorant monotone de Krasnoselskii [2], p. 401.

**PROPOSITION 4.** Soit  $E(\tau)$  un espace de Fréchet,  $K \subset E$  un cône convexe, fermé, localement borné, normal et saillant,  $f: K \rightarrow K$  un opérateur non-linéaire complètement continu.

On suppose qu'il existe un opérateur  $\mathcal{L}: K \rightarrow K$  monotone croissant et homogène tel que:

1°)  $f(x) \equiv \mathcal{L}(x); \forall x \in K$

2°)  $(\exists u_0 \in K \setminus \{0\})(\exists n \in \mathbb{N})(\exists \mu > 0)(\mathcal{L}^n(u_0) \equiv \mu u_0)$ .



Alors le couple  $(f, I)$  a une branche continue de longueur maximale de vecteurs propres sur le cône  $\mathbf{K}$ .

DÉMONSTRATION. Pour chaque  $m \in \mathbf{N}$  on considère l'opérateur complètement continu,

$$f_m(x) = f(x) + \frac{u_0}{m}; \quad \forall x \in \mathbf{K}.$$

Soit  $U$  un voisinage ouvert de zéro tel que  $U \cap \mathbf{K}$  est borné.

Le cône  $\mathbf{K}$  étant localement borné et normal, il existe [8] une seminorme continue  $p$  uniformément positive et monotone croissante sur  $\mathbf{K}$ .

Puisque la famille de semi-normes considérée sur  $E$  est suffisante et  $u_0 \in \mathbf{K} \setminus \{0\}$  il résulte que  $p(u_0) > 0$ .

Si  $x \in \partial U \cap \mathbf{K}$ , alors,

$$p(f_m(x)) \cong \inf_{y \cong m^{-1}u_0} p(y) \cong m^{-1}p(u_0) > 0.$$

Utilisant le théorème 1 il résulte qu'il existe  $x_m \in \partial U \cap \mathbf{K}$  et  $\lambda_m > 0$  tels que,

$$f(x_m) + \frac{u_0}{m} = \lambda_m x_m.$$

Comme  $\partial U \cap \mathbf{K}$  est borné et  $f$  complètement continu, on peut supposer (éventuellement considérant des sous-suites) que  $\{f(x_m)\} \rightarrow y$  et  $\{\lambda_m\} \rightarrow \lambda_*$  (quand  $m \rightarrow \infty$ ). (On observe que  $\{\lambda_m\}$  est bornée.)

On prouve que  $\{x_m\} \rightarrow x_* \neq 0$  et  $\lambda_* \neq 0$ .

En effet on a,

$$\mathcal{L}(x_m) + \frac{1}{m} u_0 \cong f(x_m) + \frac{1}{m} u_0 = \lambda_m x_m$$

et donc,  $\mathcal{L}(x_m) \cong \lambda_m x_m$  ce qui implique,

$$(*) : \mathcal{L}^n(x_m) \cong \lambda_m^n x_m.$$

De la relation,  $\frac{1}{m} u_0 \cong \lambda_m x_m$  on obtient,  $\frac{1}{m \lambda_m} u_0 \cong x_m$ , et on peut considérer le plus grand nombre  $t_m > 0$  tel que,  $t_m u_0 \cong x_m$ .

On a dans ce cas,  $t_m \mathcal{L}(u_0) \cong \mathcal{L}(x_m)$  et donc,  $t_m \mathcal{L}^n(u_0) \cong \mathcal{L}^n(x_m)$ .

Il résulte que,  $t_m \mu u_0 \cong \mathcal{L}^n(x_m) \cong \lambda_m^n x_m$  d'où,  $x_m \cong t_m \lambda_m^{-n} \mu u_0$  et utilisant la maximalité du  $t_m$  on obtient,  $t_m \lambda_m^{-n} \mu \leq t_m$ .

On a donc  $\lambda_m \cong \mu^{1/n} > 0$ ;  $\forall m \in \mathbf{N}$  d'où  $\lambda_* = \lim \lambda_m > 0$ . Comme on a,  $x_m = \lambda_m^{-1} \left[ f(x_m) + \frac{u_0}{m} \right] \xrightarrow{m \rightarrow \infty} \lambda_*^{-1} [y + 0]$  il résulte que  $\{x_m\}_{m \in \mathbf{N}}$  est convergente et si  $x_* = \lim_{m \rightarrow \infty} x_m$  on a que  $x_* \in \partial U \cap \mathbf{K}$ .

En passant à la limite dans la relation  $f(x_m) + \frac{1}{m} u_0 = \lambda_m x_m$  on obtient,  $f(x_*) = \lambda_* x_*$  où  $\lambda_* > 0$  et  $x_* \in \partial U \cap \mathbf{K}$  ce qui prouve la proposition.  $\square$



REMARQUE. Si  $u_0 \in K \setminus \{0\}$  et  $K_{u_0, k}$  est le cône considéré dans l'ouvrage [1] c'est-à-dire :

$$K_{u_0, k} = \{0\} \cup \{x \in K \mid \exists \alpha, \beta > 0, \beta/\alpha \leq k, k \geq 1, \text{ fixé } \alpha u_0 \leq x \leq \beta u_0\}$$

alors la proposition 4 s'applique dans le cas où l'opérateur  $f$  est minoré par un opérateur  $\mathcal{L}: K \rightarrow K$  tel que, il existe  $n \in \mathbb{N}$  ayant la propriété que  $\mathcal{L}^n(K) \subset K_{u_0, k}$ .

Des opérateurs ayant cette propriété sont étudiés dans l'ouvrage [1].

On termine ce paragraphe en remarquant que d'autres résultats sur l'existence des branches continues de longueur maximale de vecteurs propres pour des couples  $(f, I)$  se trouvent dans l'ouvrage [13].

Au moins du point de vue de nos résultats, l'étude de l'existence des branches continues de longueur maximale de vecteurs propres pour des couples  $(f, g)$  où  $g \neq I$ , il est très important.

5. On considère dans ce paragraphe que l'espace  $E(\tau)$  est un espace de Fréchet.

On dit que l'ensemble  $K \subset E$  est un *cône tronqué* s'il existe un cône convexe  $K_1 \subset E$  et un voisinage convexe  $V$  de zéro tel que  $K = K_1 \cap V$ .

Si, dans ce cas le cône  $K_1$  est localement borné alors on dit que l'ensemble  $K$  est un *cône tronqué localement borné*.

Soit  $E(\tau_1), F(\tau_2)$  deux espaces de Fréchet, et  $K_1 \subset E, K_2 \subset F$  deux cônes convexes ou deux cônes tronqués.

DÉFINITION 2. On dit que l'opérateur  $g: K_1 \rightarrow K_2$  est *strictement positif* si, quel que soit  $\{x_n\}_{n \in \mathbb{N}} \subset K_1$  une suite telle que  $\lim_{n \rightarrow \infty} g(x_n) = 0$  il résulte qu'il existe une sous-suite  $\{x_{n_k}\}$  convergente vers zéro.

REMARQUE. La notion d'opérateur strictement positif qu'on utilise dans cet ouvrage est plus générale que la notion d'opérateur strictement positif définie par Schaefer [28].

EXEMPLES. 1°) Les opérateurs de type  $c$ -norm considérés dans l'ouvrage [20] sont strictement positifs. Aussi dans l'ouvrage [20] il existe plusieurs exemples d'opérateurs intégraux qui sont de type  $c$ -norm. Un cas particulier d'opérateur de type  $c$ -norm est l'opérateur uniformément positif au sens de Krasnoselskii [12] c'est-à-dire un opérateur  $g: K \rightarrow K$ , ( $K$  étant un cône convexe dans un espace de Banach) tel que, il existe  $b \in \mathbb{R}_+ \setminus \{0\}$  tel que,  $(\forall x \in K)(\|g(x)\| \geq b \|x\|)$ .

2°) Soit  $E(\tau_1), F(\tau_2)$  deux espaces de Fréchet,  $K_1 \subset E$  un cône convexe et  $K_2 \subset F$  un cône convexe normal.

a) Si, pour l'opérateur  $g: K_1 \rightarrow K_2$  il existe un opérateur strictement positif  $h: K_1 \rightarrow K_2$  tel que,  $(\forall x \in K_1)(g(x) \geq h(x))$  alors  $g$  est strictement positif.

b) Si, le cône  $K_2$  est localement borné et  $p$  une semi-norme continue uniformément positive sur  $K_2$  et il existe une fonction continue  $\varphi: K_2 \rightarrow \mathbb{R}_+$  et  $a, b \in \mathbb{R}_+ \setminus \{0\}$  tels que,  $\varphi(0) = 0$  et  $(\forall x \in K_1)(\varphi(g(x)) \geq b[p(x)]^a)$ , alors  $g$  est strictement positif.

3°) Dans l'ouvrage [28] il existe plusieurs exemples d'opérateurs strictement positifs. Soit  $E(\tau_1), F(\tau_2)$  deux espaces localement convexes,  $K_1 \subset E, K_2 \subset F$  deux cônes convexes ou deux cônes tronqués. On dit que l'opérateur  $g: K_1 \rightarrow K_2$  est *propre* si et seulement si,  $g$  est continu et quel que soit  $C \subset K_2$  un ensemble compact,

l'ensemble  $g^{-1}(C)$  est compact. Si  $E$  et  $F$  sont deux espaces de Banach et  $g: E \rightarrow F$  un opérateur continu on dit que  $g$  est *coercif* si,  $\|x\| \rightarrow \infty$  implique  $\|g(x)\| \rightarrow \infty$ .

EXEMPLES. 1°) Soit  $E, F$  deux espaces de Banach et  $f: E \rightarrow F$  un opérateur continu. Si,  $(\forall x \in E)(\|f(x)\| \cong C\|x\|^a)$  où  $a, C \in \mathbb{R}_+ \setminus \{0\}$ , alors  $f$  est coercif. Les opérateurs coercifs ont été utilisés dans [31].

2°) L'opérateur de Von Kármán [2] est propre.

3°) Dans l'ouvrage [2] on trouve des exemples d'opérateurs différentiels propres.

4°) Soit  $E, F$  deux espaces de Banach et  $f: E \rightarrow F$  un opérateur coercif. Si  $f = h + g$  où  $h$  est un opérateur compact et  $g$  est un opérateur propre, alors  $f$  est propre.

5°) Soit  $E, F$  deux espaces de Banach,  $E$  étant réflexif et  $f: E \rightarrow F$  un opérateur continu et coercif. Si, quel que soit  $\{x_n\}_{n \in \mathbb{N}} \subset E$  une suite faiblement convergente vers  $x_0$  telle que  $\{f(x_n)\}_{n \in \mathbb{N}}$  est fortement convergente, il résulte que  $\{x_n\}_{n \in \mathbb{N}}$  est fortement convergente vers  $x_0$ , alors  $f$  est un opérateur propre.

DÉFINITION 3. Si  $E(\tau_1), F(\tau_2)$  sont deux espaces localement convexes et  $f: E \rightarrow F$  un opérateur continu on dit que  $f$  est *(b)-propre* si pour tout ensemble compact  $C \subset F$  et pour tout ensemble borné et fermé  $B \subset E$  l'ensemble  $f^{-1}(C) \cap B$  est compact (s'il est non-vidé).

REMARQUE. On a la même définition si  $f: K_1 \rightarrow K_2$  où  $K_1 \subset E, K_2 \subset F$  sont deux cônes convexes ou deux cônes tronqués.

EXEMPLES. 1°) Soit  $E$  un espace de Banach,  $E^*$  son dual et  $f: E \rightarrow E^*$ . On dit que  $f$  vérifie la condition  $(S_+)$  si: pour toute suite  $\{x_n\} \subset E$  faiblement convergente vers  $x$  telle que  $\overline{\lim} \langle f(x_n), x_n - x \rangle \leq 0$  il résulte que  $\{x_n\}$  est fortement convergente vers  $x$  et  $\{f(x_n)\}$  est faiblement convergente vers  $f(x)$ . Si  $E$  est un espace de Banach réflexif,  $g: E^* \rightarrow E$  un opérateur linéaire monotone (au sens de Minty—Browder),  $f: E \rightarrow E^*$  un opérateur borné qui vérifie la condition  $(S_+)$ , alors  $I + g \circ f$  est un opérateur *(b)-propre*.

2°) Chaque opérateur propre est *(b)-propre*.

3°) Soit  $E$  un espace de Banach réflexif,  $K_1 \subset E, K_2 \in E^*$  deux cônes convexes fermés. Si  $f: K_1 \rightarrow K_2$  est un opérateur continu qui vérifie la condition  $(S_0)$  alors  $f$  est *(b)-propre*. En effet, soit  $C \subset K_2$  un ensemble compact, et  $B \subset K_1$  un ensemble borné et fermé tel que  $f^{-1}(C) \cap B = \emptyset$ . Soit  $\{x_n\}_{n \in \mathbb{N}} \subset f^{-1}(C) \cap B$  une suite; éventuellement considérant une sous-suite, comme  $E$  est réflexif, on peut considérer que  $\{x_n\}_{n \in \mathbb{N}} \rightarrow x_0 \in K_1$  et  $\{f(x_n)\}_{n \in \mathbb{N}} \rightarrow y \in C$ . On peut prouver que  $\langle f(x_n), x_n \rangle \rightarrow \langle y, x_0 \rangle$  et alors la condition  $(S_0)$  implique,  $\{x_n\}_{n \in \mathbb{N}} \rightarrow x_0$ . Donc  $x_0 \in B$  et d'après la continuité de  $f$  on a,  $f(x_0) = y$  c'est-à-dire  $x_0 \in f^{-1}(C) \cap B$ .

REMARQUE. Les conditions  $(S)$  et  $(S_0)$  sont vérifiées par des opérateurs quasi-elliptiques dans des hypothèses très faibles.

6. On présente dans ce paragraphe les principaux résultats de cet ouvrage.

Les démonstrations sont données dans le cas des cônes convexes mais les résultats et les démonstrations sont vraies dans le cas des cônes tronqués.

THÉORÈME 2. Soit  $E(\tau_1), F(\tau_2)$  deux espaces de Fréchet,  $K_1 \subset E, K_2 \subset F$  deux

cônes convexes localement bornés, fermés,  $f: K_1 \rightarrow K_2$  un opérateur complètement continu et  $g: K_1 \rightarrow K_2$  un opérateur continu, strictement positif et (b)-propre.

Soit  $U_1, U_2$  deux voisinages ouverts de zéro dans l'espace  $E$  tels que  $\bar{U}_2 \subset U_1$  et  $U_1 \cap K_1$  est borné.

On suppose que les hypothèses suivantes sont vérifiées.

1°)  $\mathcal{S}(f, g)$  est une branche continue de longueur maximale de vecteurs propres sur l'ensemble  $(\bar{U}_1 \setminus U_2) \cap K_1$ .

2°) Il existe  $\mu^* > 0$  tel que:

$$i) \quad (x \in \mathcal{S}(f, g) \cap \partial U_1) \ \& \ (f(x) = \mu g(x)) \Rightarrow \mu < \mu^*$$

$$ii) \quad (x \in \mathcal{S}(f, g) \cap \partial U_2) \ \& \ (f(x) = \mu g(x)) \Rightarrow \mu > \mu^*.$$

Alors  $\mu^*$  est une valeur propre pour le couple  $(f, g)$  associée à un vecteur propre  $x^* \in K_1 \cap (U_1 \setminus \bar{U}_2)$ , c'est-à-dire,  $f(x^*) = \mu^* g(x^*)$ .

DÉMONSTRATION. Puisque le cône  $K_1$  est localement borné et  $U_1 \cap K_1$  est borné, il existe une semi-norme continue  $p$ , uniformément positive sur le cône  $K_1$  et un nombre réel  $r_0 > 0$  tel que:

$$U_1 \cap K_1 \subset \bar{B}_{r_0} = \{x \in E \mid p(x) \leq r_0\}.$$

On suppose que  $\mu^*$  n'est pas une valeur propre pour le couple  $(f, g)$  associée à un vecteur propre  $x^* \in K_1 \cap (U_1 \setminus \bar{U}_2)$  et on considère les ensembles suivants:

$$[\mathcal{S}(f, g)]^0 = \{x \in \mathcal{S}(f, g) \mid x \in \bar{U}_1 \setminus U_2\};$$

$$[\mathcal{S}(f, g)]^< = \{x \in [\mathcal{S}(f, g)]^0 \mid f(x) = \mu g(x), \mu < \mu^*\};$$

$$[\mathcal{S}(f, g)]^> = \{x \in [\mathcal{S}(f, g)]^0 \mid f(x) = \mu g(x), \mu > \mu^*\}.$$

Les hypothèses 2°) impliquent les relations suivantes:

$$[\mathcal{S}(f, g)]^0 = [\mathcal{S}(f, g)]^> \cup [\mathcal{S}(f, g)]^<;$$

$$[\mathcal{S}(f, g)]^> \cap [\mathcal{S}(f, g)]^< = \emptyset.$$

Puisque,  $[\mathcal{S}(f, g)]^>, [\mathcal{S}(f, g)]^< \subset [\mathcal{S}(f, g)]^0$  il résulte qu'il existe  $r > 0$  tel que:

$$(1): \quad \inf \{p(x) \mid x \in [\mathcal{S}(f, g)]^<\} > r \quad \text{et}$$

$$(2): \quad \inf \{p(x) \mid x \in [\mathcal{S}(f, g)]^>\} > r.$$

En effet, si on suppose que:  $\inf \{p(x) \mid x \in [\mathcal{S}(f, g)]^<\} = 0$ , alors il existe une suite  $\{x_n\}_{n \in \mathbb{N}}, \forall n \in \mathbb{N}, x_n \in [\mathcal{S}(f, g)]^<$  telle que  $\lim_{n \rightarrow \infty} p(x_n) = 0$ .

Puisque la semi-norme  $p$  est uniformément positive sur  $K$  il résulte que  $\{x_n\}_{n \in \mathbb{N}} \xrightarrow{p} 0$ .

Comme le cône  $K_1$  est fermé, il résulte que  $K_1 \cap (\bar{U}_1 \setminus U_2)$  est un ensemble fermé et comme  $[\mathcal{S}(f, g)]^< \subset K_1 \cap (\bar{U}_1 \setminus U_2)$  il résulte que  $0 \in K_1 \cap (\bar{U}_1 \setminus U_2)$  ce qui est absurde. Par un calcul analogue on prouve que,  $\inf \{p(x) \mid x \in [\mathcal{S}(f, g)]^>\} > 0$ . Donc on peut choisir  $r > 0$  qui vérifie les relations 1°), 2°). On prouve que l'ensemble  $[\mathcal{S}(f, g)]^<$  est fermé. En effet, soit  $\{x_n\}_{n \in \mathbb{N}} \subset [\mathcal{S}(f, g)]^<$  telle que:

$$(3): f(x_n) = \mu_n g(x_n) \text{ où: } 0 \leq \mu_n < \mu^* \text{ et } \{x_n\}_{n \in \mathbb{N}}$$

convergente. Soit  $x_0 = \lim_{n \rightarrow \infty} x_n$ ; puisque  $(\forall n \in \mathbb{N})(p(x_n) > r)$  il résulte que,  $\lim_{n \rightarrow \infty} p(x_n) = p(x_0) \geq r$ .

Le cône  $K_2$  étant localement borné il existe une semi-norme  $q$  continue et uniformément positive sur  $K_2$ .

On prouve maintenant la relation:

$$(4): \inf \{q(g(x)) | x \in K_1, p(x) \geq r\} = \varepsilon > 0.$$

En effet, si on suppose que,  $\inf \{q(g(x)) | x \in K_1, p(x) \geq r\} = 0$  il résulte qu'il existe une suite  $\{y_n\}_{n \in \mathbb{N}} \subset K_1$  telle que,  $(\forall n \in \mathbb{N})(p(y_n) \geq r)$  et  $\{q(g(y_n))\} \rightarrow 0$ .

La semi-norme  $q$  étant uniformément positive sur  $K_2$  il résulte que  $\{g(y_n)\}_{n \in \mathbb{N}} \xrightarrow{1} 0$ .

L'opérateur  $g$  étant strictement positif sur  $K_1$  il résulte qu'il existe une sous-suite  $\{y_{n_k}\} \xrightarrow{1} 0$  ce qui est impossible parce que,  $(\forall k \in \mathbb{N})(p(y_{n_k}) \geq r)$ .

Donc la relation (4) est vraie. De la relation (4) il résulte,

$$(5): \lim_{n \rightarrow \infty} q(g(x_n)) = q(g(x_0)) \geq \varepsilon > 0.$$

De la relation (3) on a,  $q(f(x_n)) = \mu_n q(g(x_n))$  d'où:  $\mu_n = \frac{q(f(x_n))}{q(g(x_n))}$ , ce qui implique qu'il existe,  $\mu_0 = \lim_{n \rightarrow \infty} \mu_n = \frac{q(f(x_0))}{q(g(x_0))}$ .

On a que  $\mu_0 < \mu^*$ , parce que si  $\mu_0 = \mu^*$ , comme  $x_0 \neq 0$  ( $p$  étant uniformément positive sur  $K_1$ ) on obtient que  $f(x_0) = \mu^* g(x_0)$  et on a ainsi que  $\mu^*$  est une valeur propre pour le couple  $(f, g)$  associée à un vecteur propre qui se trouve dans l'ensemble  $K_1 \cap (U_1 \setminus U_2)$  ce qui est impossible. Donc l'ensemble  $[\mathcal{S}(f, g)]^<$  est fermé.

Un calcul analogue donne aussi que l'ensemble  $[\mathcal{S}(f, g)]^>$  est fermé.

Puisque l'espace  $E(\tau_1)$  est métrisable, la topologie  $\tau_1$  est définie par une famille suffisante dénombrable de semi-normes  $\{\| \cdot \|_n\}_{n \in \mathbb{N}}$ . On peut supposer que la famille  $\{\| \cdot \|_n\}_{n \in \mathbb{N}}$  est monotone croissante et dans ce cas la topologie  $\tau_1$  est définie par la métrique:

$$\varrho(x, y) = \sum_{n=1}^{\infty} 2^{-n} \frac{|x - y|_n}{1 + |x - y|_n}; \quad \forall x, y \in E.$$

Utilisant la distance  $\varrho$  on peut définir le nombre,  $\text{dist}([\mathcal{S}(f, g)]^<, [\mathcal{S}(f, g)]^>)$ .

On démontre la relation suivante:

$$(6): \text{dist}([\mathcal{S}(f, g)]^<, [\mathcal{S}(f, g)]^>) > 0.$$

Pour prouver la relation (6) on suppose que,  $\text{dist}([\mathcal{S}(f, g)]^<, [\mathcal{S}(f, g)]^>) = 0$ .

Il existe alors,  $\{x_n\}_{n \in \mathbb{N}} \subset [\mathcal{S}(f, g)]^<$  et  $\{y_n\}_{n \in \mathbb{N}} \subset [\mathcal{S}(f, g)]^>$  telles que:  $f(y_n) = \mu_n g(y_n)$  où  $\mu_n > \mu^*$  pour tout  $n \in \mathbb{N}$  et

$$(7): \lim_{n \rightarrow \infty} \varrho(x_n, y_n) = 0. \text{ Mais on a les relations suivantes:}$$

$$(8): r < p(x_n); r < p(y_n) \leq r_0; \quad \forall n \in \mathbb{N}.$$

Puisque  $f$  est complètement continu et les relations (8) impliquent que  $\{y_n\}_{n \in \mathbb{N}}$

est bornée (à cause du fait que  $p$  est uniformément positive sur  $\mathbf{K}_1$ ) il résulte qu'il existe  $M > 0$  tel que,  $(\forall n \in \mathbf{N})(q(f(y_n)) \leq M)$ .

Les relations (4) et (8) impliquent,  $(\forall n \in \mathbf{N})(q(g(y_n)) \geq \varepsilon > 0)$  d'où il résulte:

$$\mu^* < \mu_n = \frac{q(f(y_n))}{q(g(y_n))} \leq \frac{M}{\varepsilon}; \quad \forall n \in \mathbf{N}.$$

Utilisant encore que  $f$  est complètement continu on peut affirmer qu'il existe une sous-suite convergente  $\{\mu_{n_j}\}_{j \in \mathbf{N}} \subset \{\mu_n\}_{n \in \mathbf{N}}$  et une sous-suite  $\{y_{n_j}\}_{j \in \mathbf{N}} \subset \{y_n\}_{n \in \mathbf{N}}$  telle que la suite  $\{f(y_{n_j})\}_{j \in \mathbf{N}}$  est convergente.

Evidemment, si  $\mu_0 = \lim \mu_{n_j}$  alors  $\mu_0 > 0$ .

De la relation,  $f(y_{n_j}) = \mu_{n_j} g(y_{n_j})$  il résulte que  $\{g(y_{n_j})\}_{j \in \mathbf{N}}$  est convergente et puisque  $g$  est un opérateur  $(b)$ -propre il résulte qu'il existe une sous-suite  $\{y_{n_{j_i}}\}$  convergente de la suite  $\{y_{n_j}\}$ ; soit  $y_0 = \lim_{j \rightarrow \infty} y_{n_{j_i}}$ ; évidemment  $y_0 \neq 0$  et on a,  $f(y_0) = \mu_0 g(y_0)$ .

L'ensemble  $[\mathcal{S}(f, g)]^>$  étant fermé il résulte que  $y_0 \in [\mathcal{S}(f, g)]^>$ . De la relation (7) il résulte:

$$\varrho(x_{n_{j_i}}, y_0) \leq \varrho(x_{n_{j_i}}, y_{n_{j_i}}) + \varrho(y_{n_{j_i}}, y_0) \rightarrow 0$$

ce qui donne que,  $y_0 = \lim_{j \rightarrow \infty} x_{n_{j_i}}$  et comme  $[\mathcal{S}(f, g)]^<$  est fermé on obtient que,  $y_0 \in [\mathcal{S}(f, g)]^> \cap [\mathcal{S}(f, g)]^<$  ce qui est absurde.

Donc si,  $d = \text{dist}([\mathcal{S}(f, g)]^<, [\mathcal{S}(f, g)]^>)$  on a,  $d > 0$ . Soit  $U_2^* = U_2 \cup U_2^>$  où:

$$U_2^> = \left\{ y \in E \mid \varrho(x, y) < \frac{d}{2}, x \in [\mathcal{S}(f, g)]^> \right\}.$$

L'ensemble  $U_2^*$  est un voisinage ouvert de zéro dans l'espace  $E$ . Si on considère le voisinage ouvert de zéro,  $U_2^* \cap U_1$  on a:  $U_2^* \cap U_1 \subset \bar{U}_1$  et l'hypothèse 1°) implique l'existence d'un vecteur propre  $x_0 \in \partial(U_2^* \cap U_1) \cap \mathbf{K}_1$  pour le couple  $(f, g)$ .

Puisque,  $\partial(U_2^* \cap U_1) \cap \mathbf{K}_1 \subset (\bar{U}_1 \setminus U_2) \cap \mathbf{K}_1$  il résulte que  $x_0 \in [\mathcal{S}(f, g)]^0$ . Mais d'après la construction de l'ensemble  $U_2^*$  il résulte que  $x_0 \notin [\mathcal{S}(f, g)]^>$  et aussi  $x_0 \notin [\mathcal{S}(f, g)]^<$  (parce que c'est en contradiction avec la définition du nombre  $d$ ).

On obtient ainsi une contradiction en regardant la définition et les propriétés de l'ensemble  $[\mathcal{S}(f, g)]^0$  et le théorème est démontré.  $\square$

REMARQUE. Si dans les hypothèses du théorème 2 on a en plus que  $\mu^* = 1$  on obtient un théorème d'existence pour l'équation de coïncidence,  $f(x) = g(x)$ ;  $x \in \mathbf{K}_1$ .

Le corollaire suivant c'est un théorème de point fixe sur des ensembles non-convexes dans un cône convexe d'un espace de Fréchet.

COROLLAIRE. Soit  $E(\tau)$  un espace de Fréchet,  $\mathbf{K} \subset E$  un cône convexe, fermé, localement borné et  $f: \mathbf{K} \rightarrow \mathbf{K}$  un opérateur complètement continu. Soit  $U_1, U_2$  deux voisinages ouverts de zéro tels que  $\bar{U}_2 \subset U_1$  et  $U_1 \cap \mathbf{K}$  est borné.

On suppose que les hypothèses suivantes sont vérifiées.

1°)  $\mathcal{S}(f, I)$  est une branche continue de longueur maximale de vecteurs propres sur l'ensemble  $(\bar{U}_1 \setminus U_2) \cap \mathbf{K}$ .

2°)  $(x \in \mathcal{S}(f, I) \cap \partial U_1) \ \& \ (f(x) = \mu x) \Rightarrow \mu < 1$ .



3°)  $(x \in \mathcal{S}(f, l) \cap \partial U_2) \ \& \ (f(x) = \mu x) \Rightarrow \mu > 1$ .

Alors il existe un point fixe  $x^*$  pour l'opérateur  $f$  tel que  $x^* \in K \cap (U_1 \setminus \bar{U}_2)$ .  $\square$

Soit  $E(\tau_1)$ ,  $F(\tau_2)$  deux espaces de Fréchet,  $K_1 \subset E$ ,  $K_2 \subset F$  deux cônes convexes fermés,  $\mathcal{C}$  un cône convexe muni d'une structure de Riesz et  $\alpha_1: \mathcal{B}_E \rightarrow \mathcal{C}$ ,  $\alpha_2: \mathcal{B}_F \rightarrow \mathcal{C}$  deux mesures de non-compacité précisées. Dans le cas des espaces de Fréchet on peut utiliser le cône  $\mathcal{C}$  décrit dans l'exemple 2°) du paragraphe (§ 3).

DÉFINITION 4. On dit que l'opérateur  $g: K_1 \rightarrow K_2$  est  $(\alpha_1, \alpha_2)$ -coercif si et seulement si:

- (1)  $g$  est continu.
- (2)  $g$  est borné.
- (3)  $(\exists c > 0)(\forall A \subset K_1, \text{ borné})(c\alpha_1(A) \leq \alpha_2(g(A)))$ .

PROPOSITION 5. Si  $g: K_1 \rightarrow K_2$  est  $(\alpha_1, \alpha_2)$ -coercif, alors quel que soit  $\{x_n\}_{n \in \mathbb{N}} \subset K_1$  une suite bornée telle que  $\{g(x_n)\}_{n \in \mathbb{N}}$  est convergente, il résulte que la suite  $\{x_n\}_{n \in \mathbb{N}}$  a une sous-suite convergente.

DÉMONSTRATION. Si  $\{x_n\}_{n \in \mathbb{N}}$  est bornée et  $\{g(x_n)\}_{n \in \mathbb{N}}$  est convergente, alors  $\alpha_2(\{g(x_n)\}) = 0$ .

Puisque  $g$  est  $(\alpha_1, \alpha_2)$ -coercif il résulte que  $\alpha(\{x_n\}) = 0$ , mais comme l'espace  $E$  est un espace de Fréchet il résulte que  $\{x_n\}_{n \in \mathbb{N}}$  est relativement compacte d'où la conclusion de la proposition.  $\square$

THÉORÈME 3. Soit  $E(\tau_1)$ ,  $F(\tau_2)$  deux espaces de Fréchet,  $K_1 \subset E$ ,  $K_2 \subset F$  deux cônes convexes localement bornés, fermés,  $f: K_1 \rightarrow K_2$  un opérateur  $(k, \alpha_1, \alpha_2)$ -Lipschitz et  $g: K_1 \rightarrow K_2$  un opérateur  $(\alpha_1, \alpha_2)$ -coercif et strictement positif.

Soit  $U_1, U_2$  deux voisinages ouverts de zéro dans l'espace  $E$  tels que  $\bar{U}_2 \subset U_1$  et  $U_1 \cap K_1$  est borné.

On suppose que les hypothèses suivantes sont vérifiées.

1°)  $\mathcal{S}(f, g)$  est une branche continue de longueur maximale de vecteurs propres sur l'ensemble  $(\bar{U}_1 \setminus U_2) \cap K_1$ .

2°) Il existe  $\mu^* > \frac{k}{c}$  tel que:

- i)  $(x \in \mathcal{S}(f, g) \cap \partial U_1) \ \& \ (f(x) = \mu g(x)) \Rightarrow \mu < \mu^*$ ;
- ii)  $(x \in \mathcal{S}(f, g) \cap \partial U_2) \ \& \ (f(x) = \mu g(x)) \Rightarrow \mu > \mu^*$ .

Alors il existe  $x^* \in K_1 \cap (U_1 \setminus \bar{U}_2)$  tel que,  $f(x^*) = \mu^* g(x^*)$ .

DÉMONSTRATION. On suppose qu'il n'existe aucun  $x^* \in K_2 \cap (U_1 \setminus \bar{U}_2)$  tel que  $f(x^*) = \mu^* g(x^*)$ .

Comme dans le théorème 2 on considère

$$[\mathcal{S}(f, g)]^0, [\mathcal{S}(f, g)]^<, [\mathcal{S}(f, g)]^> \text{ et on a,}$$

$$[\mathcal{S}(f, g)]^0 = [\mathcal{S}(f, g)]^> \cup [\mathcal{S}(f, g)]^< ,$$

$$[\mathcal{S}(f, g)]^> \cap [\mathcal{S}(f, g)]^< = \emptyset.$$



Puisque le cône  $K_1$  est localement borné et  $U_1 \cap K_1$  est borné, il existe une semi-norme continue  $p$  uniformément positive sur le cône  $K_1$  et un nombre réel  $r_0 > 0$  tel que:  $U_1 \cap K_1 \subset \bar{B}_{r_0} = \{x \in E \mid p(x) \leq r_0\}$ , et comme  $g$  est strictement positif le même calcul utilisé dans la démonstration du théorème 2 donne que les ensembles  $[\mathcal{S}(f, g)]^<, [\mathcal{S}(f, g)]^>$  sont fermés.

Soit  $\{\mu_n\}_{n \in \mathbb{N}}$  une famille dénombrable monotone croissante et suffisante de semi-normes qui donne la topologie  $\tau_1$  sur l'espace  $E$ .

Si on considère la métrique:

$$\varrho(x, y) = \sum_{n=1}^{\infty} 2^{-n} \frac{|x-y|_n}{1+|x-y|_n}; \quad \forall x, y \in E$$

on peut démontrer que,

$$(1): \text{dist}([\mathcal{S}(f, g)]^<, [\mathcal{S}(f, g)]^>) > 0.$$

En effet, on suppose que:  $\text{dist}([\mathcal{S}(f, g)]^<, [\mathcal{S}(f, g)]^>) = 0$ . Il existe alors  $\{x_n\}_{n \in \mathbb{N}} \subset [\mathcal{S}(f, g)]^<$  et  $\{y_n\}_{n \in \mathbb{N}} \subset [\mathcal{S}(f, g)]^>$  telles que:  $f(y_n) = \mu_n g(y_n)$ , où  $\mu_n > \mu^*$  pour tout  $n \in \mathbb{N}$  et

$$(2): \lim_{n \rightarrow \infty} \varrho(x_n, y_n) = 0.$$

Mais, comme dans le théorème 2 on peut prouver qu'il existe  $r > 0$  tel que:

$$(3): r < p(x_n); r < p(y_n) \leq r_0; \quad \forall n \in \mathbb{N}.$$

Puisque  $p$  est uniformément positive sur  $K$ , de la relation (3) il résulte que  $\{y_n\}_{n \in \mathbb{N}}$  est bornée et comme  $f$  est  $(k, \alpha_1, \alpha_2)$ -Lipschitz il résulte que  $\{f(y_n)\}_{n \in \mathbb{N}}$  est bornée.

Le cône  $K_2$  étant localement borné il existe une semi-norme  $q$  continue et uniformément positive sur  $K_2$ .

Il existe donc  $M > 0$  tel que:  $(\forall n \in \mathbb{N})(q(f(y_n)) \leq M)$ . L'opérateur  $g$  étant strictement positif on peut montrer que:

$$(4): \inf\{q(g(x)) \mid x \in K, p(x) \geq r\} = \varepsilon > 0.$$

Les relations (3) et (4) impliquent;  $(\forall n \in \mathbb{N})(q(g(y_n)) \geq \varepsilon) > 0$ , d'où il résulte la relation:

$$(5): (\forall n \in \mathbb{N}) \left( \mu^* < \mu_n = \frac{q(f(y_n))}{q(g(y_n))} \leq \frac{M}{\varepsilon} \right).$$

Donc de la relation (5) on obtient que la suite  $\{\mu_n\}_{n \in \mathbb{N}}$  est relativement compacte et alors il existe une sous-suite  $\{\mu_{n_j}\}_{j \in \mathbb{N}} \subset \{\mu_n\}_{n \in \mathbb{N}}$  convergente.

Si  $\mu_0 = \lim_{j \rightarrow \infty} \mu_{n_j}$  on a:

$$(6): \mu_0 \geq \mu^* > \frac{k}{c}.$$

Puisque  $\{f(y_{n_j})\}_{j \in \mathbb{N}}$  est bornée et  $\mu_0 = \lim_{j \rightarrow \infty} \mu_{n_j}$  on a:

$$\{(\mu_0^{-1} - \mu_{n_j}^{-1})f(y_{n_j})\} \xrightarrow{\tau_1} 0 \quad \text{ce qui implique:}$$

$$\alpha_2(\{(\mu_0^{-1} - \mu_{n_j}^{-1})f(y_{n_j})\}) = 0.$$

L'opérateur  $g$  étant  $(\alpha_1, \alpha_2)$ -coercif on a :

$$(7): \quad c\alpha_1(\{y_{n_j}\}) \equiv \alpha_2(\{g(y_{n_j})\}) \equiv \alpha_2(\{\mu_0^{-1}f(y_{n_j})\}) + \alpha_2(\{(\mu_0^{-1} - \mu_{n_j}^{-1})f(y_{n_j})\})$$

d'où :

$$c\alpha_1(\{y_{n_j}\}) \equiv \mu_0^{-1}k\alpha_1(\{y_{n_j}\}).$$

Utilisant la relation (6) on obtient alors que  $\alpha_1(\{y_{n_j}\})=0$  et comme l'espace  $E(\tau_1)$  est un espace de Fréchet il résulte que  $\{y_{n_j}\}_{j \in \mathbb{N}}$  est relativement compact. Donc il existe une sous-suite convergente de la suite  $\{y_{n_j}\}_{j \in \mathbb{N}}$ .

Nous pouvons supposer que cette sous-suite est même  $\{y_{n_j}\}_{j \in \mathbb{N}}$ . Soit  $y_0 = \lim_{j \rightarrow \infty} y_{n_j}$ ; puisque  $p(y_{n_j}) > r$  on a que  $p(y_0) \geq r$  et donc  $y_0 \neq 0$ .

Evidemment,  $g(y_0) \neq 0$  et la relation,  $f(y_{n_j}) = \mu_{n_j}g(y_{n_j})$  implique,  $f(y_0) = \mu_0g(y_0)$  et comme  $[\mathcal{S}(f, g)]^>$  est fermé on obtient,  $y_0 \in [\mathcal{S}(f, g)]^>$ .

De la relation (2) on a :

$$\varrho(x_{n_j}, y_0) \leq \varrho(x_{n_j}, y_{n_j}) + \varrho(y_{n_j}, y_0) \rightarrow 0 \quad \text{ce qui donne: } \lim_{j \rightarrow \infty} x_{n_j} = y_0$$

et comme  $[\mathcal{S}(f, g)]^<$  est fermé il résulte que  $y_0 \in [\mathcal{S}(f, g)]^> \cap [\mathcal{S}(f, g)]^<$ .

Donc, si  $d = \text{dist}([\mathcal{S}(f, g)]^>, [\mathcal{S}(f, g)]^<) > 0$  on peut considérer l'ensemble:  $U_2^* = U_2 \cup U_2^>$  où :

$$U_2^> = \left\{ y \in E \mid \varrho(x, y) < \frac{d}{2}, x \in [\mathcal{S}(f, g)]^> \right\}.$$

Comme dans le théorème 2, utilisant l'hypothèse 1°) il résulte l'existence d'un vecteur propre  $x_0 \in \mathbf{K}_1 \cap \partial(U_2^* \cap U_1)$  pour le couple  $(f, g)$  qui a les propriétés suivantes:  $x_0 \in [\mathcal{S}(f, g)]^0$ ,  $x_0 \notin [\mathcal{S}(f, g)]^>$  et  $x_0 \notin [\mathcal{S}(f, g)]^<$  ce qui est impossible et le théorème est démontré.  $\square$

**COROLLAIRE 1.** Si les hypothèses du théorème 3 sont vérifiées mais  $\mu^* = 1$  et  $\frac{k}{c} < 1$ , alors il existe  $x^* \in \mathbf{K}_1 \cap (U_1 \setminus \bar{U}_2)$  tel que  $f(x^*) = g(x^*)$ .  $\square$

Si  $E = F$  est un espace de Fréchet,  $\alpha: \mathcal{B}_E \rightarrow \mathcal{C}$  une mesure de noncompacité et  $g = I$  alors  $g$  est strictement positif sur  $\mathbf{K} \subset E$  et  $(\alpha, \alpha)$ -coercif où  $c = 1$ .

On a alors le résultat suivant.

**COROLLAIRE 2.** Soit  $E(\tau)$  un espace de Fréchet,  $\mathbf{K} \subset E$  un cône convexe, fermé, localement borné et  $f: \mathbf{K} \rightarrow \mathbf{K}$  un opérateur  $(k, \alpha, \alpha)$ -Lipschitz.

Soit  $U_1, U_2$  deux voisinages ouverts de zéro tels que,  $\bar{U}_2 \subset U_1$  et  $U_1 \cap \mathbf{K}$  est borné.

On suppose que les hypothèses suivantes sont vérifiées.

1°)  $\mathcal{S}(f, I)$  est une branche continue de longueur maximale de vecteurs propres sur l'ensemble  $(\bar{U}_1 \setminus U_2) \cap \mathbf{K}$ .

2°) Il existe  $\mu^* > k$  tel que:

$$\text{i) } (x \in \mathcal{S}(f, g) \cap \partial U_1) \ \& \ (f(x) = \mu x) \Rightarrow \mu < \mu^*;$$

$$\text{ii) } (x \in \mathcal{S}(f, g) \cap \partial U_2) \ \& \ (f(x) = \mu x) \Rightarrow \mu > \mu^*.$$

Alors,

- a) il existe  $x^* \in K \cap (U_1 \setminus \bar{U}_2)$  tel que  $f(x^*) = \mu^* x^*$ ,  
 b) si  $\mu^* = 1$  alors il existe un point fixe de l'opérateur  $f$  dans l'ensemble  $K \cap (U_1 \setminus \bar{U}_2)$ .  $\square$

**COROLLAIRE 3.** Soit  $E(\tau_1), F(\tau_2)$  deux espaces de Fréchet,  $K_1 \subset E$ ,  $K_2 \subset F$  deux cônes convexes localement bornés, fermés,  $f: K_1 \rightarrow K_2$  un opérateur  $(k, \alpha_1, \alpha_2)$ -Lipschitz et  $g: K_1 \rightarrow K_2$  un opérateur  $(\alpha_1, \alpha_2)$ -coercif et strictement positif.

Soit  $U_1, U_2$  deux voisinages ouverts de zéro dans l'espace  $E$  tels que  $\bar{U}_2 \subset U_1$  et  $U_1 \cap K_1$  est borné.

On suppose que les hypothèses suivantes sont vérifiées.

1°)  $\mathcal{S}(f, g)$  est une branche continue de longueur maximale de vecteurs propres sur l'ensemble  $(\bar{U}_1 \setminus U_2) \cap K_1$ .

$$2^\circ) \quad \gamma = \max \left( \alpha_0, \frac{k}{c} \right) < \beta \quad \text{où:}$$

$$\alpha_0 = \sup \{ \mu \mid (\exists x \in \mathcal{S}(f, g) \cap \partial U_1) \ \& \ (f(x) = \mu g(x)) \},$$

$$\beta = \inf \{ \mu \mid (\exists x \in \mathcal{S}(f, g) \cap \partial U_2) \ \& \ (f(x) = \mu g(x)) \}.$$

Alors, pour tout  $\mu^* \in ]\gamma, \beta[$  il existe  $x^* \in K_1 \cap (U_1 \setminus \bar{U}_2)$  tel que  $f(x^*) = \mu^* g(x^*)$ .  $\square$

**REMARQUES.** 1°) Les résultats de ce paragraphe, en particulier sont vraies sur des cônes convexes dans des espaces de Banach. Comme les cônes  $K_1, K_2$  ne sont pas supposés saillants il résulte que dans le cas des espaces de Banach,  $K_1$  ou  $K_2$  ou les deux peuvent être des sousespaces ou l'espace tout entier.

2°) Les résultats obtenus sont vraies sur des cônes tronqués fermés.

3°) Les théorèmes 2 et 3 peuvent être considérés comme des extensions aux couples  $(f, g)$  généraux des résultats démontrés par K. Schmitt dans l'ouvrage [30] et du théorème 5 de l'ouvrage [16] démontré par N. V. Marčenko.

4°) Dans l'ouvrage [10] nous avons donné une signification intuitive des théorèmes 2 et 3 dans le cas  $E = F$ .

5°) Si on analyse la démonstration du théorème 2 (resp. du théorème 3) on constate que le théorème 2 (resp. le théorème 3) reste vraie si on demande que la branche  $\mathcal{S}(f, g)$  a seulement la propriété suivante:  $\partial U \cap \mathcal{S}(f, g) \cap A \neq \emptyset$  seulement pour chaque voisinage ouvert  $U$  de zéro tel que, 1)  $U \supset U_2$ , 2)  $U \cap K_1$  est borné, 3)  $\partial U \subset A$ .

Cette observation est importante parce que, il est possible d'avoir que le couple  $(f, g)$  a une branche continue de vecteurs propres sur le cône  $K_1$  sans que la définition 1 soit vérifiée sur  $A$ , en temps que les conditions précédentes sont vérifiées.

7. On présente maintenant quelques applications des théorèmes 2 et 3.

a) Soit  $f: [0, +\infty[ \rightarrow \mathbf{R}$  une fonction continue telle que,  $f(a) > 0$  et  $f(b) < 0$  où:  $a, b \in \mathbf{R}$ ,  $0 < a < b$ , alors on sait qu'il existe  $x_* \in ]a, b[$  tel que  $f(x_*) = 0$ . On montre maintenant que ce résultat classique est un corollaire du théorème 2. Soit  $m_f = \inf \{ f(x) \mid x \in [a, b] \}$  et on considère la fonction  $g(x) = mx$ ;  $\forall x \in [0, +\infty[$  où  $m$  est choisi tel que,  $g(x) > -m_f$ ;  $\forall x \in [a, b]$ . On observe que  $g$  est continue,

propre et strictement positive au sens de la définition 2. On a aussi les relations suivantes:

$$a_1) \quad f(x) + g(x) > 0; \quad \forall x \in [a, b],$$

$$a_2) \quad g(x) > 0; \quad \forall x \in [a, b],$$

$$a_3) \quad \text{l'opérateur } f+g: \mathbf{R}_+ \rightarrow \mathbf{R} \text{ est complètement continue.}$$

On considère les voisinages de zéro,  $U_1 = ]-\varepsilon, b[$ ,  $U_2 = ]-\varepsilon, a[$ ;  $\varepsilon > 0$  et le couple d'opérateurs  $(f+g, g)$ . Comme pour chaque  $\alpha \in [a, b]$  il existe  $\lambda > 0$  tel que  $f(\alpha) + g(\alpha) = \lambda g(\alpha)$  il résulte que  $\mathcal{S}(f+g, g)$  est une branche continue de vecteurs propres de longueur maximale sur  $(U_1 \setminus U_2) \cap \mathbf{R}_+$ . Soit  $\mu^* = 1 > 0$ ; alors on a,

$$f(b) + g(b) = \mu g(b) \rightarrow \mu < 1,$$

$$f(a) + g(a) = \mu g(a) \rightarrow \mu > 1.$$

Donc, toutes les hypothèses du théorème 2 sont vérifiées d'où il résulte qu'il existe  $x_* \in ]a, b[$  tel que,  $f(x_*) + g(x_*) = g(x_*)$ , c'est-à-dire  $f(x_*) = 0$ .  $\square$

b°) De l'application précédente il résulte que le théorème 2 peut être utilisé pour construire un algorithme par bisection pour approximer les solutions de l'équation  $f(x) = 0$  dans un espace de Banach  $E$ , ordonné par un cône convexe  $\mathbf{K}$  et où,  $f: \mathbf{K} \rightarrow F$ ,  $F$  étant un espace de Banach aussi. On suppose que par une localisation grossière il existe  $U_1, U_2$  deux voisinages ouverts de zéro tels que,  $\bar{U}_2 \subset \subset U_1$ ,  $U_1 \cap \mathbf{K}$  est borné et que l'équation  $f(x) = 0$  a une solution dans l'ensemble  $(U_1 \setminus \bar{U}_2) \cap \mathbf{K}$ . On suppose aussi qu'il existe un opérateur  $g: \mathbf{K} \rightarrow F$  tel que  $(f+g, g)$  vérifie les hypothèses du théorème 2 pour  $U_1$  et  $U_2$ . On divise l'ensemble  $(U_1 \setminus \bar{U}_2) \cap \mathbf{K}$  par un autre voisinage ouvert de zéro  $U_3$  tel que:  $\bar{U}_2 \subset U_3$ ,  $\bar{U}_3 \subset U_1$ . Si pour le couple  $(f+g, g)$  et  $\mu^* = 1$  les hypothèses 2°) du théorème 2 sont vérifiées pour  $U_2$  et  $U_3$  alors la solution de l'équation  $f(x) = 0$  se trouve dans l'ensemble  $(U_3 \setminus \bar{U}_2) \cap \mathbf{K}$ , si non, elle se trouve dans l'ensemble  $U_1 \setminus U_3$ . On continue ainsi d'appliquer le théorème 2 à l'ensemble retenu jusqu'à ce qu'on obtient une localisation très fine de la solution étudiée.  $\square$

c°) Les dernières années plusieurs auteurs [2], [6], [18], [22], [23], [5] ont consacré leurs ouvrages à l'étude des équations non linéaires qui sont en réalité des cas particuliers de l'équation: (\*):  $f(x) = \lambda g(\lambda, x)$  où  $f: E \rightarrow F$ ,  $g: \mathbf{R} \times E \rightarrow F$ ;  $E$  et  $F$  étant deux espaces de Banach. L'équation (\*) peut être étudiée utilisant le théorème 2 (ou le théorème 3). En effet, on considère l'espace  $\tilde{E} = \mathbf{R} \times E$ , le cône convexe  $\mathbf{K} = \mathbf{R}_+ \times E$  et les opérateurs:  $\tilde{f}: \mathbf{K} \rightarrow F$ ,  $\tilde{g}: \mathbf{K} \rightarrow F$  définis par:

$$\tilde{f}(\lambda, x) = f(x); \quad \forall (\lambda, x) \in \mathbf{K},$$

$$\tilde{g}(\lambda, x) = \lambda g(\lambda, x); \quad \forall (\lambda, x) \in \mathbf{K}.$$

Si les hypothèses du théorème 2 (ou du théorème 3) sont vérifiées pour le couple  $(\tilde{f}, \tilde{g})$ ,  $\mu^* = 1$  et pour les voisinages de zéro  $U_1$  et  $U_2$  bien choisis dans l'espace  $\tilde{E}$  alors il résulte qu'il existe  $(\lambda^*, x^*) \in (U_1 \setminus \bar{U}_2) \cap \mathbf{K}$  tel que:  $f(x^*) = \lambda^* g(\lambda^*, x^*)$ .  $\square$

d°) Un modèle important utilisé en économie mathématique pour étudier les phénomènes de croissance est le modèle de J. von Neumann qui conduit à un pro-

blème de valeurs propres de la forme  $Ax \leq \lambda Bx$ ,  $A$  et  $B$  étant deux matrices. Si le phénomène de croissance soit-il économique ou biologique n'est pas linéaire donc perturbé alors on a un problème de valeurs propres de la forme,  $f(x) \leq \lambda g(x)$ . Ce problème peut être étudié par comparaison utilisant le théorème 2. En effet, soit  $E, F$  deux espaces de Banach ordonnés par  $K_1 \subset E$  (resp.  $K_2 \subset F$ ) deux cônes convexes saillants, et  $f, g: K_1 \rightarrow K_2$  deux opérateurs non-linéaires. Par certaines estimations on soupçonne qu'il est possible que  $\lambda^* > 0$  soit un facteur de croissance, c'est-à-dire qu'il existe un vecteur de croissance  $x_* \in K_1 \setminus \{0\}$  tel que  $f(x_*) \leq \lambda^* g(x_*)$ . Pour prouver ça on utilise le théorème 2. On suppose qu'il existe  $f_1, g_1$  tels que:  $f(x) \leq f_1(x); g_1(x) \leq g(x); \forall x \in K_1$ . Soit  $U_1, U_2$  deux voisinages de zéro comme dans le théorème 2. Si les hypothèses du théorème 2 sont vérifiées pour  $f_1, g_1, U_1, U_2$  et  $\lambda^*$ , alors il existe  $x_* \in (U_1 \setminus \bar{U}_2) \cap K_1$  tel que  $f(x_*) \leq f_1(x_*) = \lambda^* g_1(x_*) \leq \lambda^* g(x_*)$ . Donc  $\lambda^*$  est un facteur de croissance et on a aussi une estimation du vecteur de croissance associé.  $\square$

## REFERENCES

- [1] BAHTIN, I. A., On the existence of eigenvectors for positive noncompletely continuous linear operators, *Math. Sb. (N. S.)* **64** (1964), 102—114 (in Russian). *MR* **29** #6286. = *Amer. Math. Soc. Translations* **51** (1966), 201—214.
- [2] BERGER, M. S., *Nonlinearity and functional analysis*, Academic Press, New York—London, 1977. *MR* **58** #7671.
- [3] DEAN, E. T. and CHAMBRE, P. L., On the bifurcation of solutions of the nonlinear eigenvalue problem  $Lu + \lambda b(x)u = g(x, u)$ , *SIAM J. Appl. Math.* **20** (1971), 722—734. *MR* **45** #3976.
- [4] FITZPATRICK, P. M. and PETRYSHYN, W. V., Fixed point theorems and the fixed point index for multivalued mappings in cones, *J. London Math. Soc.* **12** (1975), 78—85. *MR* **53** #8974.
- [5] FITZPATRICK, P. M. and PETRYSHYN, W. V., Positive eigenvalues for nonlinear multivalued noncompact operators with applications to differential operators, *J. Differential Equations* **22** (1976), 428—441. *MR* **55** #8909.
- [6] FITZPATRICK, P. M. and PETRYSHYN, W. V., On the nonlinear eigenvalue problem  $Tu = \lambda Cu$  involving noncompact abstract and differential operators, *Boll. Un. Mat. Ital. B* (5) **15** (1978), 80—107. *MR* **80b**: 47083.
- [7] GAINES, R. E. and MAWHIN, J. L., *Coincidence degree and nonlinear differential equations*, Lecture Notes in Math., Springer-Verlag, Berlin, 1977. *MR* **58** #30551.
- [8] ISAC, G., *Cônes localement bornés, théorèmes de point fixe et valeurs propres positives pour des opérateurs non-linéaires dans des espaces localement convexes*, Inst. de Math. Pure et Appl., Univ. Catholique de Louvain. Rapport Nr. 131 (1979).
- [9] ISAC, G., Valeurs propres positives pour des opérateurs de type  $\alpha$ -Lipschitz dans des espaces de Fréchet, *Publ. Math. de l'Univ. de Pau* (1979).
- [10] ISAC, G., Equations de coïncidence sur des cônes et branches continues de vecteurs propres, *Seminari dell'Istituto di Matematica Applicata „Giovanni Sansone” Univ. degli studi di Firenze*, Facoltà d'ingegneria (1981).
- [11] JAMESON, G. O., *Ordered linear spaces*, Lecture Notes in Math., Bd. 141, Springer-Verlag, Berlin, 1970. *MR* **55** #10996.
- [12] KRASNOSELSKII, M. A., *Positive solutions of operator equations*, Noordhoff, Groningen, 1964. *MR* **31** #6107.
- [13] KRASNOSELSKII, M. A., *Topological methods in the theory of nonlinear integral equations*, Pergamon, Oxford, 1964. *MR* **28** #2414.
- [14] KRASNOSELSKII, M. A., Some problems of nonlinear analysis, *Uspehi Mat. Nauk.* **3** (61) (1954), 57—114. *MR* **17**—769.
- [15] MANGASARIAN, O. L., Perron—Frobenius of  $Ax - \lambda Bx$ , *J. Math. Anal. Appl.* **36** (1971), 86—102. *MR* **44** #2773.
- [16] MARČENKO, N. V., On the continuation of an operator and the existence of fixed points, *Dokl. Akad. Nauk* **145** (1962), 1767—1770.
- [17] MAWHIN, J. L., Topological degree methods in nonlinear boundary value problems, NSF—



- CBMS Regional Conference at the Claremont University Center, Claremont, California (1977).
- [18] MCCORMICK, S. F., A mesh refinement method for  $Ax = \lambda Bx$ , *Math. Comp.* **36** (1981), 485—498. *MR* **82d**: 65070.
  - [19] MININNI, M., Coincidence degree and solvability of some nonlinear functional equations in normed spaces. A spectral approach, *Nonlinear Anal.* **1** (1977), 105—122. *MR* **58** # 31193.
  - [20] NÉMETH, A. B., Nonlinear operators that transform a wedge, Babeş—Bolyai University, Faculty of Mathematics, Research Seminars, Preprint Nr. 1—1980, Cluj-Romania, 27—43.
  - [21] PERESSINI, A. L., *Ordered topological vector spaces*, Harper & Row, New York—London, 1967. *MR* **37** # 3315.
  - [22] PETRYSHYN, W. V., Bifurcation and asymptotic bifurcation for equations involving  $A$ -proper mappings with applications to differential equations, *J. Differential Equations* **28** (1978), 124—154. *MR* **57** # 51426.
  - [23] PETRYSHYN, W. V., On the solvability of  $x \in Tx + \lambda Fx$  in quasinormal cones with  $T$  and  $F$   $k$ -set-contractive, *Nonlinear Anal.* **5** (1981), 585—591. *MR* **82f**: 47074.
  - [24] PETRYSHYN, W. V., Nonlinear eigenvalue problems and the existence of nonzero fixed points for  $A$ -proper mappings. (In: *Theory of Nonlinear Operators*, Proc. of Intern. Summer School at Berlin, G.D.R., Ed. R. Kluge, Akad. Verlag, Berlin, 1978, 215—227). *MR* **80d**: 47002.
  - [25] PETRYSHYN, W. V. and FITZPATRICK, P. M., A degree theory, fixed point theorems and mapping theorems for multivalued non-compact mappings, *Trans. Amer. Math. Soc.* **194** (1974), 1—25.
  - [26] SADOVSKI, B. N., Limit-compact and condensing operators, *Uspehi Mat. Nauk* **27** (1972), 81—146. *MR* **55** # 1161.
  - [27] SCOTT, D. S., Solving sparse symmetric generalized eigenvalue problems without factorization, *Siam J. Numer. Anal.* **18** (1981), 102—110.
  - [28] SCHAEFER, H. H., On non-linear positive operators, *Pacific J. Math.* **9** (1959), 847—860. *MR* **22** # 1827.
  - [29] SCHECHTER, M., SHAPIRO, J. and SNOW, M., Solution of the nonlinear problem  $A(u) = N(u)$  in Banach space, *Trans. Amer. Math. Soc.* **241** (1978), 69—78. *MR* **81g**: 47069.
  - [30] SCHMITT, K., Fixed point and coincidence theorems with applications to nonlinear differential and integral equations, *Inst. de Math. Pure et Appl. Univ. Catholique de Louvain, Rapport Nr. 97* (1976).
  - [31] FUČIK, S., NEČAS, J., SOUČEK, J. and SOUČEK, V., *Spectral analysis of nonlinear operators*, Lecture Notes in Math., Bd. 396, Springer-Verlag, Berlin, 1973. *MR* **57** # 7280.

(Recu le 3 janvier 1983)

DÉPARTEMENT DE MATHÉMATIQUES  
COLLÈGE MILITAIRE ROYAL  
SAINT-JEAN, QUÉBEC  
JOUÏRO  
CANADA



# SOME BIVARIATE EXTENSIONS OF THE GENERALIZED WARING DISTRIBUTION

EVDOKIA XEKALAKI

## Summary

The bivariate generalized Waring distribution results as a mixture of the double Poisson distribution. In this paper some probability models are considered that give rise to alternative bivariate forms of the generalized Waring distribution.

## 1. Introduction

The bivariate generalized Waring distribution with parameters  $a, k, m$  and  $\varrho$  (B.G.W.D.  $(a; k, m; \varrho)$ ) defined by Xekalaki [7] is the distribution with probability function (p. f.)  $p_{r,l}$  given by

$$(1.1) \quad p_{rl} = \frac{\varrho(k+m)}{(a+\varrho)_{(k+m)}} \frac{a_{(r+l)} k_{(r)} m_{(l)}}{(a+k+\varrho+m)_{(r+l)}} \frac{1}{r!} \frac{1}{l!}$$

$$r = 0, 1, 2, \dots, \quad l = 0, 1, 2, \dots$$

where  $a, k, m, \varrho > 0$  and  $\alpha_{(\beta)} = \Gamma(\alpha + \beta)/\Gamma(\alpha)$ ,  $\alpha > 0$ ,  $\beta \in \mathbb{R}$ . The probability generating function (p.g.f.) of this distribution is

$$G(s, t) = \frac{\varrho(k+m)}{(a+\varrho)_{(k+m)}} F_1(a; k, m; a+k+\varrho+m; s, t)$$

where  $F_1(a; b, b'; x, y)$  is the Appell function of the first kind defined by

$$F_1(a; b, b'; c; x, y) = \sum_{m,n} \frac{a_{(m+n)} b_{(m)} b'_{(n)}}{c_{(m+n)}} \frac{x^m}{m!} \frac{y^n}{n!}$$

$$a, b, b', c - a - b - b' > 0, \quad (x, y) \in [-1, 1] \times [-1, 1].$$

The marginal probability distributions of  $X$  and  $Y$ , the conditional distributions of  $X|Y=y$  and  $Y|X=x$  as well as the distribution of  $X+Y$  are all of the same form. Specifically, they are univariate generalized Waring distributions (U.G.W.D.) with p.g.f.'s expressed in terms of the Gauss hypergeometric function obtained from

$$(1.2) \quad {}_pF_q(a_1, a_2, \dots, a_p; b_1, b_2, \dots, b_q; z) = \sum_{r=0}^{\infty} \frac{(a_1)_{(r)} \dots (a_p)_{(r)}}{(b_1)_{(r)} \dots (b_q)_{(r)}} \frac{z^r}{r!}$$

1980 *Mathematics Subject Classification*. Primary 60E05; Secondary 62H05.

*Key words and phrases*. Bivariate generalized Waring distributions, Gauss hypergeometric function, confluent hypergeometric function.

for  $p=2$ ,  $q=1$ . The U.G.W.D. is a member of the family of generalized hypergeometric distributions studied by, among others, Kemp and Kemp [2] and Sarkadi [3]. (For more information concerning the structure, properties and applications of the U.G.W.D.  $(a, k; q)$  the interested reader is referred to Irwin [1] and Xekalaki [4], [5], [6], [8].)

Xekalaki [7] showed that the B.G.W.D. can arise as the joint distribution of accidents incurred in two consecutive time periods by a group of people in situations where not only random factors are present, but, also, factors associated with the individual's exposure to external risk as well as psychological factors predisposing the individual to accidents. In the same paper, except for providing a satisfactory fit to accident data, the B.G.W.D. was shown to enable one to separately estimate the variance components due to random, external and psychological factors so that one can have a clue as to which kind of factors influenced the particular accident situation the most.

The derivation of the B.G.W.D. in the context of the above-mentioned accident situation was based on a mixed Poisson model. The mathematical nature of the mixing processes involved suggests the possibility of obtaining some more general forms of bivariate distributions with both marginals of the U.G.W.D. type. This can be done by slightly altering the mixing mechanisms concerned.

In the sequel, we provide various mixed models which give rise to such distributions. Some properties and limiting cases of these models are also considered.

## 2. Some bivariate extensions of the generalized Waring distribution

We first give a description of Xekalaki's [7] model.

MODEL 1 (Xekalaki, [7]). Let  $X_1$ ,  $X_2$  be nonnegative integer valued random variables (r.v.'s) whose joint distribution is the double Poisson with p.g.f.

$$(2.1) \quad g(s, t) = e^{\lambda_1(s-1) + \lambda_2(t-1)}, \quad \lambda_1, \lambda_2 > 0.$$

Assume that  $\lambda_1$  and  $\lambda_2$  are independent gamma r.v.'s with probability density functions

$$(2.2) \quad f(\lambda_1) = \frac{v^{-k}}{\Gamma(k)} e^{-(1/v)\lambda_1} \lambda_1^{k-1}, \quad v, k > 0$$

and

$$(2.3) \quad g(\lambda_2) = \frac{v^{-m}}{\Gamma(m)} e^{-(1/v)\lambda_2} \lambda_2^{m-1}, \quad m > 0,$$

respectively. Then (2.1) becomes

$$(2.4) \quad \begin{aligned} G(s, t) &= \frac{v^{-k}}{\Gamma(k)} \int_0^\infty e^{-(\lambda_1/v)(1+v(1-s))} \lambda_1^{k-1} d\lambda_1 \frac{v^{-m}}{\Gamma(m)} \int_0^\infty e^{-(\lambda_2/v)(1+v(1-t))} \lambda_2^{m-1} d\lambda_2 = \\ &= (1+v(1-s))^{-k} (1+v(1-t))^{-m}, \end{aligned}$$

i.e.  $(X_1, X_2)|v$  follows a double negative binomial distribution with parameters  $k, v/(1+v), m$  and  $v/(1+v)$ . Assume now that  $v$  has a beta distribution of the second kind (beta II) with parameters  $a$  and  $q$  and p.d.f.

$$h(v) = \frac{\Gamma(a+q)}{\Gamma(q)\Gamma(a)} v^{a-1} (1+v)^{-(a+q)}, \quad a, q > 0.$$

Then the resulting distribution of  $X_1, X_2$  has p.g.f.

$$\begin{aligned} G_{x_1 x_2}(s, t) &= \frac{\Gamma(q+a)}{\Gamma(q)\Gamma(a)} \int_0^\infty v^{a-1} (1+v)^{-(a+q)} (1+v(1-s))^{-k} (1+v(1-t))^{-m} dv = \\ &= \frac{q(k+m)}{(a+q)_{(k+m)}} F_1(a; k, m; a+k+m+q; s, t), \end{aligned}$$

i.e.,  $(X_1, X_2) \sim \text{B.G.W.D.}(a; k, m; q)$ .

Xekalaki [8] has shown that under certain conditions the B.G.W.D.  $(a; k, m; q)$  tends to the double negative binomial distribution with parameters  $k, m, \frac{a}{a+q}$  and  $\frac{a}{a+q}$ , and to the double Poisson distribution with parameters  $\frac{ak}{a+q}$  and  $\frac{am}{a+q}$ . Moreover, if the scale is at our choice, the B.G.W.D. can be shown to tend to the bivariate beta II distribution with parameters  $k, m$  and  $q$  or to an uncorrelated bivariate gamma distribution with parameters  $k, m, 1$  and  $1$ .

MODEL 2. Let  $X_1, X_2$  be non-negative discrete r.v.'s whose joint distribution is the double Poisson with parameters  $\lambda p$ , and  $\lambda q$  and p.g.f. given by

$$(2.5) \quad g(s, t) = e^{\lambda(p(s-1)+q(t-1))}, \quad \lambda, p, q > 0, \quad p+q \leq 1.$$

Assume that  $\lambda$  is a r.v. having a gamma distribution with p.d.f.

$$(2.6) \quad f(\lambda) = \frac{(a/b)^a}{\Gamma(a)} e^{-(a/b)\lambda} \lambda^{a-1}, \quad \lambda, a, b > 0.$$

Then (2.5) becomes

$$\begin{aligned} (2.7) \quad G(s, t) &= \frac{(a/b)^a}{\Gamma(a)} \int_0^\infty e^{-\lambda(a/b)(1+(b/a)p(1-s)+(b/a)q(1-t))} \lambda^{a-1} d\lambda = \\ &= \left[ 1 + \frac{b}{a} \{p(1-s) + q(1-t)\} \right]^{-a}, \end{aligned}$$

i.e.,  $(X_1, X_2)|b$  follows a bivariate negative binomial distribution with parameters  $a, bp/(a+bp)$  and  $bq/(a+bq)$ .

If we now let  $b$  have a Beta II distribution with parameters  $k, q$  and p.d.f.

$$\begin{aligned} (2.8) \quad h(b) &= \frac{\Gamma(q+k)a^{-1}}{\Gamma(q)\Gamma(k)} \left(\frac{b}{a}\right)^{k-1} \left(1+\frac{b}{a}\right)^{-(q+k)} \\ &\quad q > 0, k > 0, b > 0, a > 0 \end{aligned}$$

the final resulting joint distribution of  $X_1, X_2$  has p.g.f.

$$\begin{aligned}
 G_{X_1, X_2}(s, t) &= \frac{\Gamma(\varrho+k)a^{-1}}{\Gamma(\varrho)\Gamma(k)} \int_0^{\infty} \left(\frac{b}{a}\right)^{k-1} \left(1+\frac{b}{a}\right)^{-(\varrho+k)} \times \\
 (2.9) \quad &\times \left[1 + \frac{b}{a} \{p(1-s) + q(1-t)\}\right]^{-a} db = \\
 &= \frac{\varrho_{(k)}}{(a+\varrho)_{(k)}} {}_2F_1(a, k; a+k+\varrho; 1-p-q+ps+qt).
 \end{aligned}$$

The generalized hypergeometric series in (2.9) is convergent for all  $(s, t) \in [-1, 1] \times [-1, 1]$ .

It can be seen that

$$(2.10) \quad G_{X_1}(s) = G_{X_1, X_2}(s, 1) = \frac{\varrho_{(k)}}{(a+\varrho)_{(k)}} {}_2F_1(a, k; a+k+\varrho; 1-p+ps)$$

with a similar expression for the p.g.f. of  $X_2$ .

$$\begin{aligned}
 G_{X_1|X_2}(s) &= \frac{\partial^{x_2}}{\partial t^{x_2}} G_{X_1, X_2}(s, 0) \bigg/ \frac{\partial^{x_2}}{\partial t^{x_2}} G_{X_1, X_2}(1, 0) = \\
 (2.11) \quad &= \frac{{}_2F_1(a+x_2, k+x_2; a+k+\varrho+x_2; 1-p-q+ps)}{{}_2F_1(a+x_2, k+x_2; a+k+\varrho+x_2; 1-q)}
 \end{aligned}$$

with a similar expression for the p.g.f. of  $X_2|X_1$  and

$$(2.12) \quad G_{X_1+X_2}(s) = G_{X_1, X_2}(s, s) = \frac{\varrho_{(k)}}{(a+\varrho)_{(k)}} {}_2F_1(a, k; a+k+\varrho; 1+(p+q)(s-1)),$$

i.e., the marginals and the convolution have all the same form.

We note that (2.10) is in fact a U.G.W.D. generalized by a binomial distribution with index 1 and probability of success  $p$  and can arise in situations where sampling is made with inclusion probability equal to  $p$ . Hence, (2.10) defines a more general distribution which includes the U.G.W.D. as a special case ( $p=1$ ). It is interesting to see that the factorial moments of the distribution generated by (2.10) are given by

$$(2.13) \quad \mu_{(r)}(x) = p^r \frac{a_{(r)}k_{(r)}}{(q-1)\dots(q-r)}.$$

It is also of interest to remark that in the case  $p+q=1$ , (2.9) and (2.12) reduce to

$$(2.14) \quad G_{X_1, X_2}(s, t) = \frac{\varrho_{(k)}}{(a+\varrho)_{(k)}} {}_2F_1(a, k; a+k+\varrho; ps+qt)$$

and

$$(2.15) \quad G_{X_1+X_2}(s) = \frac{\varrho_{(k)}}{(a+\varrho)_{(k)}} {}_2F_1(a, k; a+k+\varrho; s) \sim \text{U.G.W.D. } (a, k; \varrho).$$

Moreover, (8.1.7) becomes

$$G_{X_1|X_2}(s) = \frac{{}_2F_1(a+x_2, k+x_2; a+k+\varrho+x_2; ps)}{{}_2F_1(a+x_2, k+x_2; a+k+\varrho+x_2; p)}$$

which, provided that  $q > x_2$ , is a weighted U.G.W.D.  $(a + x_2, k + x_2; q - x_2)$  and can arise in cases where the sampling chance (weight) is proportional to  $p^{x_1}$ ,  $p \leq 1$ .

In the more general case when  $p + q \leq 1$  it can be shown, using an argument similar to that used by Xekalaki [8] that

(i) for  $q \rightarrow \infty$ ,  $a \rightarrow \infty$ ,  $\frac{a}{a+q} < 1$ ,  $0 < k < \infty$ , the distribution defined by (2.9)

tends to the bivariate negative binomial distribution with parameters  $k$ ,  $\frac{ap}{a+q}$  and  $\frac{aq}{a+q}$ .

(ii) if we let  $a \rightarrow \infty$ ,  $k \rightarrow \infty$ ,  $q \rightarrow \infty$  while  $\frac{a}{a+q} \rightarrow 0$  and  $\frac{ka}{a+q} < \infty$ , we obtain

the double Poisson distribution with parameters  $\frac{akp}{a+q}$  and  $\frac{akq}{a+q}$ .

MODEL 3. Consider  $X_1$  and  $X_2$  to be two non-negative discrete r.v.'s whose joint distribution is the double Poisson with parameters  $\lambda_1, \lambda_2$  and p.g.f. given by

$$(2.16) \quad g(s, t) = e^{\lambda_1(s-1) + \lambda_2(t-1)}, \quad \lambda_1, \lambda_2 \geq 0.$$

Assume that  $(\lambda_1, \lambda_2)$  is a random vector having a bivariate gamma distribution with p.d.f.

$$(2.17) \quad f(\lambda_1, \lambda_2) = \frac{b^{-m} e^{-(1/b)\lambda_2}}{\Gamma(k)\Gamma(m-k)} \lambda_1^{k-1} (\lambda_2 - \lambda_1)^{m-k-1}$$

$$b > 0, m > k > 0, \lambda_2 \geq \lambda_1 > 0.$$

Then, from (2.16) we have

$$(2.18) \quad G(s, t) = \frac{b^{-m}}{\Gamma(k)\Gamma(m-k)} \int_0^\infty \int_0^{\lambda_2} e^{-(1/b)\lambda_2(1+b(1-t) + \lambda_1(1-s))} \times \\ \times \lambda_1^{k-1} (\lambda_2 - \lambda_1)^{m-k-1} d\lambda_1 d\lambda_2 = \\ = \frac{1}{\Gamma(k)\Gamma(m-k)} \int_0^\infty e^{-(\lambda_2/b)(1+b(1-t))} \left(\frac{\lambda_2}{b}\right)^{m-k-1} \int_0^{\lambda_2/b} e^{\lambda_1(1-s)} \times \\ \times \left(\frac{\lambda_1}{b}\right)^{k-1} \left(1 - \frac{\lambda_1}{\lambda_2}\right)^{m-k-1} d\frac{\lambda_1}{b} d\frac{\lambda_2}{b} = \frac{1}{\Gamma(k)\Gamma(m-k)} \int_0^\infty e^{-(\lambda_2/b)(1+b(1-t))} \left(\frac{\lambda_2}{b}\right)^{m-1} \times \\ \times \int_0^1 e^{-\lambda_1 \lambda_2(1-s)} \lambda_1^{k-1} (1 - \lambda_1)^{m-k-1} d\lambda_1 d\frac{\lambda_2}{b} = \\ = \frac{1}{\Gamma(m)} \int_0^\infty {}_1F_1(k; m; \lambda_2(s-1)) e^{-(\lambda_2/b)(1+b(1-t))} \left(\frac{\lambda_2}{b}\right)^{m-1} d\frac{\lambda_2}{b} = \\ = [1+b(1-t)]^{-m} {}_1F_0\left(k; ; \frac{b(s-1)}{1+b(1-t)}\right) = \\ = [1+b(1-t)]^{-m+k} [1+b(1-t)+b(1-s)]^{-k}$$

where  ${}_1F_1(a; b; z)$  and  ${}_1F_0(a; ; z)$  are obtained from (1.2) for  $p=q=1$  and  $p=1, q=0$ , respectively.

If we now consider  $b \sim \text{Beta II}(a; \varrho)$ , i.e.

$$(2.19) \quad b \sim \frac{\Gamma(a+\varrho)}{\Gamma(a)\Gamma(\varrho)} b^{a-1}(1+b)^{-(a+\varrho)}, \quad b > 0, a, \varrho > 0,$$

we obtain the final joint distribution of  $X_1, X_2$  as

$$(2.20) \quad G_{X_1 X_2}(s, t) = \frac{\Gamma(a+\varrho)}{\Gamma(a)\Gamma(\varrho)} \int_0^\infty b^{a-1}(1+b)^{-(a+\varrho)} [1+b(1-t)]^{-m+k} \times \\ \times [1+b(1-t)+b(1-s)]^{-k} db = \frac{\varrho_{(m)}}{(a+\varrho)_{(m)}} F_1(a; m-k, k; a+\varrho+m; t, t+s-1).$$

Note that the  $F_1$  series in (2.20) is convergent for  $(s, t) \in [-1, 1] \times [-1, 1]$ .

We have for the marginal distributions of  $X_1$  and  $X_2$

$$(2.21) \quad G_{X_1}(s) = \frac{\varrho_{(m)}}{(a+\varrho)_{(m)}} F_1(a; m-k, k; a+\varrho+m; 1, s) =$$

$$= \frac{\varrho_{(k)}}{(a+\varrho)_{(k)}} {}_2F_1(a, k; a+\varrho+k; s) \sim \text{U.G.W.D.}(a, k; \varrho)$$

and

$$(2.22) \quad G_{X_2}(t) = \frac{\varrho_{(m)}}{(a+\varrho)_{(m)}} F_1(a; m-k, k; a+\varrho+m; t, t) =$$

$$= \frac{\varrho_{(m)}}{(a+\varrho)_{(m)}} {}_2F_1(a, m; a+\varrho+m; t) \sim \text{U.G.W.D.}(a, m; \varrho),$$

respectively.

Thus, both marginal distributions of (2.20) are U.G.W.D.'s. This was expected due to the fact that the distributions of  $\lambda_1$  and  $\lambda_2$  are, from (2.17), gamma  $(k, b^{-1})$  and gamma  $(m, b^{-1})$ , respectively.

The conditional distributions of  $X_2|X_1, X_1|X_2$  and the distribution of  $X_1+X_2$ , however, are not of a U.G.W.D. form. Nor are they of a more general form containing the U.G.W.D. as a particular case.

It is interesting to see that for  $a \rightarrow \infty, \varrho \rightarrow \infty$  while  $\frac{a}{a+\varrho} \rightarrow q < 1, 0 < k, m < \infty$  the distribution in (2.20) tends to a distribution with p.g.f.

$$\{1+q(1-t)/p\}^{-m+k} \{1+q[(1-t)+(1-s)]/p\}^{-k}, \quad p = 1-q.$$



Moreover, for  $a \rightarrow \infty, k \rightarrow \infty, m \rightarrow \infty, \varrho \rightarrow \infty$  while  $\frac{a}{a+\varrho} \rightarrow 0, \frac{ak}{a+\varrho} < \infty$  and  $\frac{am}{a+\varrho} < \infty$  it reduces to the double Poisson distribution with parameters  $\frac{ak}{a+\varrho}$  and  $\frac{am}{a+\varrho}$ .

MODEL 4. Let  $X_1, X_2$  be non-negative discrete r.v.'s and let their joint distribution be the double Poisson with parameters  $\lambda_1, \lambda_2$  and p.d.f. given by (2.16). Assume that

$$(2.23) \quad \lambda_1 \sim \frac{\Gamma(\varrho+k)b^{-1}}{\Gamma(\varrho)\Gamma(k)} \left(\frac{\lambda_1}{b}\right)^{k-1} \left(1+\frac{\lambda_1}{b}\right)^{-(\varrho+k)}, \quad \lambda_1 > 0, \quad \varrho, k, b > 0$$

$$(2.24) \quad \lambda_2 \sim \frac{\Gamma(\varrho+m)b^{-1}}{\Gamma(\varrho)\Gamma(m)} \left(\frac{\lambda_2}{b}\right)^{m-1} \left(1+\frac{\lambda_2}{b}\right)^{-(\varrho+m)}, \quad \lambda_2 > 0, \quad \varrho, m, b > 0.$$

Then (2.16) becomes

$$(2.25) \quad \begin{aligned} G(s, t) &= \frac{\Gamma(\varrho+k)\Gamma(\varrho+m)}{\{\Gamma(\varrho)\}^2\Gamma(k)\Gamma(m)} \int_0^\infty e^{\lambda_1(s-1)} \left(\frac{\lambda_1}{b}\right)^{k-1} \left(1+\frac{\lambda_1}{b}\right)^{-(\varrho+k)} d\frac{\lambda_1}{b} \times \\ &\times \int_0^\infty e^{\lambda_2(t-1)} \left(\frac{\lambda_2}{b}\right)^{m-1} \left(1+\frac{\lambda_2}{b}\right)^{-(\varrho+m)} d\frac{\lambda_2}{b} = \\ &= \frac{\Gamma(\varrho+k)\Gamma(\varrho+m)}{\{\Gamma(\varrho)\}^2} \psi(k; 1-\varrho; (1-s)b) \psi(m; 1-\varrho; (1-t)b) \end{aligned}$$

where  $\psi$  is the confluent hypergeometric function of the second kind defined by

$$\psi(a; c; z) = \frac{1}{\Gamma(a)} \int_0^\infty t^{a-1} (1+t)^{c-a-1} e^{-zt} dt, \quad a > 0, \quad c-a > 0.$$

Letting  $b$  be a r.v. with p.d.f.

$$(2.26) \quad f(b) = \frac{1}{\Gamma(a)} e^{-b} b^{a-1}, \quad a > 0, \quad b > 0$$

we obtain the p.g.f. of the final resulting distribution of  $(X_1, X_2)$  as

$$\begin{aligned} G_{X_1, X_2}(s, t) &= \frac{\Gamma(\varrho+k)\Gamma(\varrho+m)}{\{\Gamma(\varrho)\}^2\Gamma(a)} \int_0^\infty e^{-b} b^{a-1} \psi(k; 1-\varrho; b(1-s)) \times \\ &\times \psi(m; 1-\varrho; b(1-t)) db. \end{aligned}$$

It can be shown that

$$\psi(a; c; z) = \frac{\Gamma(1-c)}{\Gamma(1+a-c)} {}_1F_1(a; c; z).$$

Then,

$$\begin{aligned}
 G_{X_1, X_2}(s, t) &= \frac{1}{\Gamma(a)} \int_0^\infty e^{-b} b^{a-1} {}_1F_1(k; 1-q; b(1-s)) \times \\
 (2.27) \quad &\quad \times {}_1F_1(m; 1-q; b(1-t)) db = \\
 &= \frac{1}{\Gamma(a)} \sum_{r,l} \frac{k_{(r)} m_{(l)}}{(1-q)_{(r)}(1-q)_{(l)}} \frac{(1-s)^r}{r!} \frac{(1-t)^l}{l!} \int_0^\infty e^{-b} b^{a+r+l-1} db = \\
 &= \sum_{r,l} \frac{a_{(r+l)} k_{(r)} m_{(l)}}{(1-q)_{(r)}(1-q)_{(l)}} \frac{(1-s)^r}{r!} \frac{(1-t)^l}{l!} = F_2(a; k, m; 1-q, 1-q; 1-s, 1-t)
 \end{aligned}$$

where  $F_2$  is as defined by

$$F_2(a; b, b'; c, c'; x, y) = \sum_{m,n} \frac{a_{(m+n)} b_{(m)} b'_{(n)}}{c_{(m)} c'_{(n)}} \frac{x^m}{m!} \frac{y^n}{n!}.$$

The region of convergence for the p.g.f. given by (2.27) is  $|1-s| + |1-t| < 1$ . Clearly, (2.16), (2.23), (2.24) and (2.26) imply that

$$(2.28) \quad X_1 \sim \text{Poisson}(\lambda_1) \sim \text{Beta II}(k; q) \sim \text{gamma}(a; 1) \sim \text{U.G.W.D.}(a, k; q),$$

$$(2.29) \quad X_2 \sim \text{Poisson}(\lambda_2) \sim \text{Beta II}(m; q) \sim \text{gamma}(a; 1) \sim \text{U.G.W.D.}(a, m; q).$$

Again, the convolution  $X_1 + X_2$  and the conditional distributions of  $X_2|X_1$ ,  $X_1|X_2$  are not U.G.W.D.'s.

We note that the double negative binomial and the double Poisson can be obtained as the limit of (2.27) for suitable limiting values of the parameters. In particular, for  $a \rightarrow \infty$ ,  $q \rightarrow \infty$  while  $\frac{a}{a+q} < 1$ ,  $0 < k < \infty$  and  $0 < m < \infty$ , the double negative binomial  $\left(k, m; \frac{a}{a+q}, \frac{a}{a+q}\right)$  distribution arises. Also, for  $a \rightarrow \infty$ ,  $k \rightarrow \infty$ ,  $m \rightarrow \infty$ ,  $q \rightarrow \infty$  while  $\frac{a}{a+q} \rightarrow 0$ ,  $\frac{ak}{a+q} < \infty$  and  $\frac{am}{a+q} < \infty$  we obtain the double Poisson  $\left(\frac{ak}{a+q}, \frac{am}{a+q}\right)$  distribution.

## REFERENCES

- [1] IRWIN, J. O., The Generalized Waring Distribution, *J. Roy. Statist. Soc. Ser. A* **138** (1975), 18—31 (Part I). *MR* **56** # 4012; 204—227 (Part II). *MR* **56** # 4013; 374—384 (Part III). *MR* **56** # 4014.
- [2] KEMP, C. D. and KEMP, A. W., Generalized hypergeometric distributions, *J. Roy. Statist. Soc. Ser. B* **18** (1956), 202—211. *MR* **18**—769.
- [3] SARKADI, K., Generalized Hypergeometric Distributions, *Magyar Tud. Akad. Mat. Kutató Int. Közl.* **2** (1957), 59—69. *MR* **20** # 6149.
- [4] XEKALAKI, E., Chance Mechanisms for the Univariate Generalized Waring Distribution and Related Characterizations, *Statistical Distributions in Scientific Work, Vol. 4* (Models, Structures and Characterizations). C. Taillie, B. Baldessari and G. P. Patil (eds.), D. Reidel Publishing Company, 1981, 157—171. *MR* **83i**: 62044.
- [5] XEKALAKI, E., The Univariate Generalized Waring Distribution in Relation to Accident Theory: Proneness, Spells or Contagion?, *Biometrics* **37** (1983), 887—895.

- [6] XEKALAKI, E., Infinite divisibility, Completeness and Regression Properties of the Univariate Generalized Waring Distribution, *Ann. Inst. Statist. Math.* **35** (1983), 161—171.
- [7] XEKALAKI, E., The Bivariate Generalized Waring Distribution and its Application to Accident Theory, *J. Roy. Statist. Soc. Ser. A* **147** (1984), 488—498.
- [8] XEKALAKI, E., Models Leading to the Bivariate Generalized Waring Distribution, *Utilitas Math.* **25** (1984), 263—290.

(Received February 22, 1983)

DEPARTMENT OF STATISTICS AND ACTUARIAL SCIENCE  
THE UNIVERSITY OF IOWA  
IOWA CITY, IA 52242  
U.S.A.

Present address :

DEPARTMENT OF STATISTICS AND INFORMATICS  
THE ATHENS SCHOOL OF ECONOMICS AND  
BUSINESS SCIENCE  
76, PATISSION STREET  
GR—104-34 ATHENS  
GREECE



# THE DISTRIBUTION OF PRIMITIVE ABUNDANT NUMBERS

ALEKSANDAR IVIĆ

*Dedicated to Prof. P. Erdős on the occasion of his seventieth birthday*

## Abstract

This note contains improvements of classical bounds for  $N(x)$ , due to P. Erdős. Here  $N(x)$  denotes the number of primitive abundant numbers not exceeding  $x$ , where  $n$  is said to be primitive abundant if  $\sigma(n) \geq 2n$  and  $\sigma(d) < 2d$  for all divisors  $d$  of  $n$  which are less than  $n$ .

## 1. Introduction

A natural number  $n$  is called primitive abundant (p.a.n. for shortness) if  $\sigma(n) \geq 2n$  and  $\sigma(d) < 2d$  for all divisor  $d$  of  $n$  which are less than  $n$ . Here as usual  $\sigma(n)$  is the sum of divisor function. In a classical paper P. Erdős [3] proved the remarkable result

$$(1.1) \quad x \exp(-8(\log x \log \log x)^{1/2}) < N(x) < x \exp\left(-\frac{1}{25}(\log x \log \log x)^{1/2}\right),$$

where  $N(x)$  denotes the number of p.a.n.'s not exceeding  $x$ . Although almost fifty years have passed since [3] was published, no sharpening of (1.1) seems to have been attained yet, and the aim of this note is to provide a proof of the following

**THEOREM.** For  $x \geq x_0(\varepsilon)$

$$(1.2) \quad x \exp(-(6^{1/2} + \varepsilon)(\log x \log \log x)^{1/2}) < N(x) < x \exp(-(12^{-1/2} - \varepsilon)(\log x \log \log x)^{1/2}).$$

Here  $\varepsilon$  denotes positive constants which may be made arbitrarily small, and which are not necessarily the same at each occurrence. The proof of (1.2) consists of refining at certain places the original proof given by Erdős [3] and using the estimate

$$(1.3) \quad \Psi(x, y) = \sum_{n \leq x, P(n) \leq y} 1 \ll x \log^2 y \exp(-u \log u + O(u \log \log u)).$$

This follows from the work of N.G. de Bruijn [1]; here  $P(n)$  is the largest prime factor of  $n$ ,  $n > 1$ ,  $3 < u = \log x / \log y < 4y^{1/2} / \log y$ ,  $2 \leq y \leq x$ . Instead of (1.3) we could use sharper estimates (e.g. see [2]), but this would not give a better upper bound

---

1980 *Mathematics Subject Classification*. Primary 10H25.

*Key words and phrases*. Primitive abundant numbers, sum of divisor functions, largest prime factor of an integer, squarefull and squarefree parts of an integer.

for  $N(x)$  than the one in (1.2). It seems reasonable to conjecture that

$$(1.4) \quad N(x) = x \exp(-(B+o(1))(\log x \log \log x)^{1/2})$$

for some  $B$  satisfying  $12^{-1/2} < B < 6^{1/2}$ , and an asymptotic formula of this sort does indeed appear for certain sums involving the function  $\Psi(x, y)$  (see [4]).

## 2. Proof of the upper bound

In what follows let  $L = C(\log x \log_2 x)^{1/2}$ ,  $\log_2 x = \log \log x$ , where  $C > 0$  will be suitably chosen later. Note that using (1.3) we obtain

$$(2.1) \quad \Psi(x, \exp L) \ll x \exp(-(1/2C - \varepsilon)(\log x \log_2 x)^{1/2}),$$

and so we may restrict ourselves to the p.a.n.'s not exceeding  $x$  for which

$$(2.2) \quad P(n) > \exp L.$$

Now every  $n \geq 2$  may be uniquely written as  $n = qs$ ,  $q = q(n)$ ,  $s = s(n)$ ,  $(q, s) = 1$ , where  $q$  is squarefree and  $s$  is squarefull (meaning  $p|s$  implies  $p^2|s$  for any prime  $p$ ). The number of squarefull  $s$  not exceeding  $x$  is  $O(x^{1/2})$ , hence the number of  $n \leq x$  for which

$$(2.3) \quad s = s(n) \leq \exp((1/3 - \varepsilon)L)$$

does not hold is by partial summation

$$(2.4) \quad \ll x \sum_{s > \exp((1/3 - \varepsilon)L)} s^{-1} \ll x \exp(-(C/6 - \varepsilon)(\log x \log_2 x)^{1/2}).$$

Thus we may suppose that  $n$  is a p.a.n. not exceeding  $x$  for which (2.2) and (2.3) hold. We shall show that the squarefree part  $q = q(n)$  of such  $n$  has a divisor in the interval

$$(2.5) \quad I = \left[ \exp(L/6), \frac{1}{2} \exp(L/2) \right].$$

The proof is by contradiction. Namely by (2.3) we have  $n = uv$ , where we may suppose  $P(u) < \exp(L/6)$  and  $p > (1/2) \exp(L/2)$  for all primes  $p$  dividing  $v$ . Denoting  $S(n) = \sigma(n)/n$ ,  $S(n)$  is multiplicative and  $S(u) = \sigma(u)/u \leq 2 - 1/u$ , since  $u$  is a divisor of a p.a.n. and consequently  $\sigma(u) \leq 2u - 1$ . If  $u \leq \exp((1/2 - \varepsilon)L)$ , then

$$\begin{aligned} 2 &\leq S(n) = S(u)S(v) \leq (2 - 1/u) \prod_{p|v} (1 + 1/p) \leq \\ &\leq (2 - \exp((\varepsilon - 1/2)L))(1 + 2 \exp(-L/2))^{\omega(v)} \leq \\ &\leq (2 - \exp((\varepsilon - 1/2)L))(1 + 4 \log x \exp(-L/2)), \end{aligned}$$

where we used the trivial bound  $\omega(v) \leq \log x$  for  $\omega(v)$ , the number of distinct prime divisors of  $v$ . But for  $x \geq x_0(\varepsilon)$  the last expression is less than 2, proving that  $u > \exp((1/2 - \varepsilon)L)$ . Let  $q(u) = p_1 p_2 \dots p_t$  ( $p_1 < \dots < p_t$  primes) denote the squarefree part of  $u$ . By (2.3)  $s = s(n)$  divides  $u$  and  $s \leq \exp((1/3 - \varepsilon)L)$ , hence we have



$q(u) \cong \exp(L/6)$ . Since for  $i=1, \dots, t$  we have  $p_i < \exp(L/6)$  by the construction of  $u$ , there must exist an  $m (\cong t-1)$  such that

$$p_1 \dots p_m \cong \exp(L/6) < p_1 \dots p_m p_{m+1}.$$

Then setting  $d = p_1 \dots p_m p_{m+1}$  we see that  $d$  is a divisor of  $q(u)$  which lies in the interval  $I$  given by (2.5), and this is a contradiction. Thus our hypothesis is not true, i.e.  $I$  does always contain a divisor of  $q(n)$  as asserted.

Having established this fact we resume the proof of the upper bound in (1.2). If (2.2) and (2.3) hold, then  $P^2(n)$  does not divide  $n$  and

$$\begin{aligned} (2.6) \quad 2 \cong S(n) &= S(P(n))S(n/P(n)) = (1 + 1/P(n))S(n/P(n)) \cong \\ &\cong 2 + 2 \exp(-L). \end{aligned}$$

Let  $C_1, C_2, \dots, C_M$  be p.a.n.'s not exceeding  $x$  which satisfy (2.2) and (2.3), and let  $D_1, D_2, \dots, D_M$  be the respective divisors of their squarefree parts lying in the interval  $I$ . Then for  $i=1, 2, \dots, M$  we have

$$(2.7) \quad C_i/D_i \ll x \exp\left(-\frac{C}{6}(\log x \log_2 x)^{1/2}\right),$$

and if we can show that all numbers  $C_i/D_i$  are distinct, then we shall obtain

$$(2.8) \quad M \ll x \exp\left(-\frac{C}{6}(\log x \log_2 x)^{1/2}\right).$$

Suppose that  $C_i/D_i = C_j/D_j$  for some  $i \neq j \leq M$ . Then

$$S(C_i/D_i) = S(C_j/D_j),$$

while

$$S(C_i) = S(C_i/D_i)S(D_i), \quad S(C_j) = S(C_j/D_j)S(D_j),$$

so that

$$(2.9) \quad S(C_i)/S(C_j) = S(D_i)/S(D_j).$$

Since  $D_i$  and  $D_j$  are squarefree and  $D_i \neq D_j$ , then it is readily checked that  $S(D_i) \neq S(D_j)$ , and we may assume  $S(D_i) > S(D_j)$ , say. Now

$$S(D_i)/S(D_j) = \sigma(D_i)D_j/(\sigma(D_j)D_i) > 1$$

implies

$$\sigma(D_i)D_j \cong \sigma(D_j)D_i + 1,$$

and since  $\sigma(D_i) < 2D_i$  we infer that

$$(2.10) \quad S(C_i)/S(C_j) = S(D_i)/S(D_j) \cong \frac{\sigma(D_j)D_i + 1}{\sigma(D_j)D_i} \cong 1 + (2D_i D_j)^{-1} > 1 + \exp(-L),$$

since  $D_i \in I$  and  $D_j \in I$ . However, from (2.6) we have

$$S(C_i)/S(C_j) \cong (2 + 2 \exp(-L))/2 = 1 + \exp(-L),$$

which contradicts (2.10) and proves  $C_i/D_i \neq C_j/D_j$ . Therefore from (2.1), (2.4) and (2.8) we obtain

$$N(x) < x \exp\left(-(1/2C - \varepsilon)(\log x \log_2 x)^{1/2}\right) + x \exp\left(-(C/6 - \varepsilon)(\log x \log_2 x)^{1/2}\right).$$

The exponential terms are of the same order of magnitude if  $1/2C = C/6$ , i.e. for  $C = 3^{1/2}$ , which yields the upper bound in (1.2).

### 3. Proof of the lower bound

To prove the lower bound in (1.2) we shall employ the ingenious construction given by Erdős in [3], but we shall choose the parameters more carefully. Namely if  $p_1 < \dots < p_k$  are any  $k$  primes between  $(k-1)2^{l+1}$  and  $k2^{l+1}$  (where of course  $l$  must be chosen in such a way that the interval in question does contain at least  $k$  primes), then it may be checked that all numbers  $m = 2^l p_1 \dots p_k$  are indeed primitive abundant, and it remains to choose suitably  $k = k(x)$  and  $l = l(x)$  and estimate from below the number of such  $m$ 's not exceeding  $x$ . To this end let

$$(3.1) \quad y = D(\log x \log_2 x)^{1/2}, \quad z = D^{-1}(\log x / \log_2 x)^{1/2},$$

where  $D > 0$  will be suitably chosen, and let

$$(3.2) \quad e^{y-1} < 2^l \leq e^y, \quad k = [z - 1 - 1/(2D^2) - \varepsilon],$$

so that  $e^{yz} = x$ , and for  $x \geq x_0(\varepsilon)$

$$z \log z < (1/(2D) + D\varepsilon/3)(\log x \log_2 x)^{1/2}.$$

Then we have

$$(3.3) \quad 2^l p_1 \dots p_k < e^y z^z e^{(y+1)(z-1-1/(2D^2)-\varepsilon)} \leq \\ \leq x \exp\{y + z \log z - (1/2D + D + D\varepsilon)(\log x \log_2 x)^{1/2} + O(z)\} \leq x \exp\left(-\frac{\varepsilon}{2}y\right) < x.$$

By the prime number theorem we obtain

$$\pi(k2^{l+1}) - \pi((k-1)2^{l+1}) = \int_{(k-1)2^{l+1}}^{k2^{l+1}} \frac{dt}{\log t} + O(k2^l \exp(-l^{1/2})) \leq \\ \leq \frac{2^{l+1}}{\log k + (l+1) \log 2} + O(2^l l^{-2}) > 2^{l+1} l^{-1}.$$

For  $r \geq k$  we have

$$\binom{r}{k} = \frac{r!}{k!(r-k)!} \geq (r/k)^k,$$

and so with  $r = [2^{l+1} l^{-1}] > e^{y-1} y^{-1}$  we obtain, for  $x \geq x_0(\varepsilon)$ ,

$$N(x) \geq \binom{r}{k} \geq \exp\{(y-1-\log y)(z-1-1/(2D^2)-\varepsilon) - z \log z + O(\log_2 x)\} \geq \\ \geq x \exp\left\{-(D + 3/(2D) + 2D\varepsilon)(\log x \log_2 x)^{1/2}\right\}.$$

The optimal choice for  $D$  is  $D^2 = 3/2$ ,  $D + 3/(2D) = 6^{1/2}$ , which gives the lower bound in (1.2).

## REFERENCES

- [1] DE BRUIJN, N. G., On the number of integers  $\leq x$  and free of prime factors  $> y$ , *Indag. Math.* **13** (1951), 50—60. *MR* **13**—724.
- [2] CANFIELD, E. R., ERDŐS, P. and POMERANCE, C., On a problem of Oppenheim concerning "Factorisatio Numerorum", *J. Number Theory* **17** (1983), 1—28. *MR* **85j**: 11012.
- [3] ERDŐS, P., On primitive abundant numbers, *J. London Math. Soc.* **9** (1935), 49—58.
- [4] Ivić, A., Sum of reciprocals of the largest prime factor of an integer, *Arch. Math. (Basel)* **36** (1981), 57—61. *MR* **82g**: 10067.

(Received April 12, 1983)

KATEDRA MATEMATIKE  
RUDARSKO-GEOLOŠKI FAKULTET  
UNIVERSITETA U BEOGRADU  
DJUŠINA 7  
YU—11000 BEOGRAD  
YUGOSLAVIA



# ON SQUAREFREE NUMBERS WITH RESTRICTED PRIME FACTORS

ALEKSANDAR IVIĆ

*Dedicated to Prof. Paul Erdős on the occasion of his seventieth birthday*

## Abstract

This note answers to a certain extent a question raised by P. Erdős concerning the distribution of squarefree integers  $n$  not exceeding  $x$ , all of whose prime factors do not exceed  $y$ . This function is being compared to the number of positive integers  $n$  not exceeding  $x$ , all of whose prime factors do not exceed  $y$ .

Let us define  $P(n)$  as the largest prime factor of  $n$ ,  $P(1)=1$ , and let

$$(1) \quad \Psi_2(x, y) = \sum_{n \leq x, P(n) \leq y} \mu^2(n)$$

denote the number of squarefree integers not exceeding  $x$ , all of whose prime factors do not exceed  $y$ . This function is naturally to be compared with the well-known function

$$(2) \quad \Psi(x, y) = \sum_{n \leq x, P(n) \leq y} 1,$$

which represents the number of integers not exceeding  $x$ , all of whose prime factors do not exceed  $y$ . Prof. P. Erdős has asked me in correspondence to investigate the validity of the asymptotic formula

$$(3) \quad \Psi_2(x, y) = (6\pi^{-2} + o(1)) \Psi(x, y), \quad (x, y \rightarrow \infty).$$

This is a natural question, since the density of squarefree numbers is  $6\pi^{-2}$ . We shall prove in this note the following result, which establishes (3) in a slightly sharper form for a suitable range of  $y$ .

**THEOREM.** *For any fixed  $\varepsilon > 0$  and*

$$(4) \quad \exp((\log \log x)^{2+\varepsilon}) \leq y \leq x$$

*we have*

$$(5) \quad \Psi_2(x, y) = \{6\pi^{-2} + O((\log \log x)^{-\varepsilon})\} \Psi(x, y).$$

Before proceeding to the proof of the Theorem we may remark that (3) cannot hold for all  $y \leq x$ , e.g. it certainly cannot hold for  $1 \ll y \leq \log x$ . To see this, observe that we have

$$(6) \quad \Psi_2(x, y) \leq \sum_{k=0}^{\pi(y)} \binom{\pi(y)}{k} = 2^{\pi(y)},$$

1980 *Mathematics Subject Classification*. Primary 10H15.

*Key words and phrases*. Squarefree numbers, largest prime factor of an integer, the Möbius function.

where by the prime number theorem

$$\pi(y) = \sum_{p \leq y} 1 = (1 + o(1)) \frac{y}{\log y}.$$

On the other hand, P. Erdős [3] proved

$$\Psi(x, \log x) = 4^{(1+o(1)) \log x / \log \log x}, \quad (x \rightarrow \infty)$$

so that  $\Psi(x, \log x)$  is roughly  $\Psi_2(x, \log x)$  squared. Also using (6) and the simple bound

$$\Psi(x, y) \cong \binom{\pi(y) + u}{u}, \quad u = [\log x / \log y],$$

(see N. G. de Bruijn [2] or P. Erdős and J. van Lint [4]) we infer that (3) cannot hold for  $1 \ll y < \log x$ .

The proof of the Theorem will employ the asymptotic formula

$$(7) \quad \Psi(x, y) = \left(1 + O_\varepsilon \left( \frac{u \log(u+1)}{\log x} \right)\right) x \varrho(u),$$

which holds uniformly in the range

$$(8) \quad x \geq 3, \quad 1 \leq u = \log x / \log y \leq \log x / (\log \log x)^{5/3 + \varepsilon}.$$

This was recently proved by A. Hildebrand [5], who improved earlier results of N.G. de Bruijn [2] and H. Maier (unpublished). The function  $\varrho(u)$  is defined as  $\varrho(u) = 1$  for  $0 \leq u \leq 1$ ,  $u \varrho'(u) = -\varrho(u-1)$  for  $u \geq 1$ , and for  $u \rightarrow \infty$  N. G. de Bruijn [1] has obtained the asymptotic formula

$$\varrho(u) = \exp \left\{ -u(\log u + \log_2 u - 1 + (\log_2 u - 1)/\log u + O((\log_2 u / \log u)^2)) \right\},$$

where  $\log_2 u = \log \log u$ . To prove the Theorem we start from

$$\mu^2(n) = \sum_{d^2 | n} \mu(d) = \begin{cases} 1 & \text{if } n \text{ is squarefree,} \\ 0 & \text{otherwise.} \end{cases}$$

Then we may write

$$\begin{aligned} \Psi_2(x, y) &= \sum_{n \leq x} \sum_{\substack{P(n) \leq y \\ d^2 | n}} \mu(d) = \\ (9) \quad &= \sum_{d \leq x^{1/2}} \sum_{\substack{P(d) \leq y}} \mu(d) \sum_{\substack{m \leq x/d^2 \\ P(m) \leq y}} 1 = \sum_{d \leq x^{1/2}} \sum_{\substack{P(d) \leq y}} \mu(d) \Psi(x/d^2, y) = \\ &= \sum_{d \leq x^{1/2}} \sum_{\substack{P(d) \leq y}} \left\{ 1 + O((\log \log x)^{-1-\varepsilon}) \right\} \mu(d) x d^{-2} \varrho(u - 2 \log d / \log y), \end{aligned}$$

where in view of (4) and (8) we used (7) to estimate  $\Psi(x/d^2, y)$ . To estimate the terms of the last sum in (9) for which  $d > D = (\log \log x)^\varepsilon$ , we shall use the inequality

$$(10) \quad \varrho(u-t) \ll \varrho(u) (u \log^2(u+1))^t.$$



This was proved by Hildebrand [5] to hold uniformly for  $u \geq 1$  and  $0 \leq t \leq u$ . Using (10) we obtain

$$\begin{aligned} \sum_{d \leq D} d^{-2} \varrho(u - 2 \log d / \log y) &\ll \varrho(u) \sum_{d \leq D} d^{-2} \exp \left\{ 2(\log u + 2 \log \log(u+1)) \frac{\log d}{\log y} \right\} \ll \\ &\ll \varrho(u) D^{-1} \exp(3D \log \log x / (\log y)) \ll \varrho(u) (\log \log x)^{-\varepsilon} \end{aligned}$$

which gives

$$\Psi_2(x, y) = \sum_{d \leq D} (1 + O(D^{-1})) x \mu(d) d^{-2} \varrho(u - 2 \log d / \log y) + O(x \varrho(u) D^{-1}).$$

Next note that  $\varrho(u)$  is non-increasing for  $u \geq 0$ , and that from its definition one has

$$\varrho(a) - \varrho(b) = \int_a^b \varrho(t-1) t^{-1} dt.$$

Thus using (10) again we have, for  $1 \leq d \leq D = (\log \log x)^c$ ,

$$\begin{aligned} \varrho(u) - \varrho(u - 2 \log d / \log y) &\ll \varrho(u - 1 - 2 \log d / \log y) \log d / (u \log y) \ll \\ &\ll \varrho(u) u^{2 \log d / \log y} (\log \log x)^{2+4 \log d / \log y} \log d / \log y \ll \varrho(u) \log d (\log \log x)^{-\varepsilon}. \end{aligned}$$

Therefore we have in view of (7)

$$\begin{aligned} \Psi_2(x, y) &= (1 + O(D^{-1})) x \sum_{d \leq D} \mu(d) d^{-2} \varrho(u) + O(x \varrho(u) D^{-1}) = \\ &= (6\pi^{-2} + O(D^{-1})) x \varrho(u) = \{6\pi^{-2} + O((\log \log x)^{-\varepsilon})\} \Psi(x, y), \end{aligned}$$

which was to be proved. A reasonable conjecture is that (5) (or at least (3)) holds in the range  $\log^{1+\varepsilon} x \leq y \leq x$ .

It may be also remarked that the foregoing proof may be easily generalized to hold for  $k$ -free numbers also. Recall that for  $k \geq 2$  fixed a number is  $k$ -free if it is not divisible by any  $p^k$ ,  $p$  prime, and

$$\sum_{d^k | n} \mu(d) = \begin{cases} 1 & \text{if } n \text{ is } k\text{-free,} \\ 0 & \text{otherwise.} \end{cases}$$

Repeating the above proof and using  $\sum_{n=1}^{\infty} \mu(n) n^{-k} = 1/\zeta(k)$ , we obtain

$$\Psi_k(x, y) = \{1/\zeta(k) + O((\log \log x)^{-\varepsilon})\} \Psi(x, y)$$

for  $y$  satisfying (4), where  $\Psi_k(x, y)$  is the number of  $k$ -free numbers not exceeding  $x$ , all of whose prime factors do not exceed  $y$ .

ACKNOWLEDGEMENT. I wish to thank the referee for useful remarks and Mat. Inst. and Rep. Zaj. of Serbia (Belgrade) for financing this research.

NOTE added in proof (25 July, 1986). Recently, the author and G. Tenenbaum (Local densities over integers free of large prime factors, *Quart. J. Math. Oxford Ser. (2)* 37 (1986), 401—417) proved that (3) holds if and only if  $\log y / \log \log x \rightarrow \infty$  as  $x \rightarrow \infty$ .

## REFERENCES

- [1] DE BRUIJN, N. G., The asymptotic behaviour of a function occurring in the theory of primes, *J. Indian Math. Soc. (N. S.)* **15** (1951), 25—32. *MR* **13**—326.
- [2] DE BRUIJN, N. G., On the number of positive integers  $\leq x$  and free of prime factors  $> y$ , I—II., *Indag. Math.* **13** (1951), 50—60. *MR* **13**—724; **28** (1966), 239—247. *MR* **34** #5770.
- [3] ERDŐS, P., *Wisk. Opgaven* **21** (1963), 133—135, problem and solution No. 136.
- [4] ERDŐS, P. and VAN LINT, J., On the number of positive integers  $\leq x$  and free of prime factors  $> y$ , *Simon Stevin* **40** (1966/67), 73—76. *MR* **35** #2836.
- [5] HILDEBRAND, A., On the number of positive integers  $\leq x$  and free of prime factors  $> y$ , *J. Number Theory* **22** (1986), 289—307.

(Received April 12, 1983)

KATEDRA MATEMATIKE  
RUDARSKO-GEOLOŠKI FAKULTET  
UNIVERSITETA U BEOGRADU  
DJUŠINA 7  
YU—11000 BEOGRAD  
YUGOSLAVIA

## STEINNESS AND THE VANISHING OF COHOMOLOGY

GUNNAR BERG

### Introduction

The main result of this note is that a relatively compact, open subspace  $X$  of a Stein space is a domain of holomorphy if the Dolbeault cohomology groups  $H^p(X, \mathcal{O})$  vanish for  $p > 0$ .

Arguably the most important result in the theory of Stein manifolds and Stein spaces is the so-called Theorem B of Cartan and Serre. This theorem states that the cohomology groups  $H^p(X, \mathcal{F})$ , where  $X$  is a Stein space and  $\mathcal{F}$  is a coherent analytic sheaf on  $X$ , vanish if  $p > 0$  (see for example [4], p. 243).

Among the many consequences we mention:

(i) Theorem A, i.e. global sections of coherent analytic sheaves generate the stalks, [7];

(ii) holomorphic functions on closed subvarieties extend to the whole space ([4], p. 245);

(iii) the first Cousin problem is always solvable ([4], p. 248).

As for the vanishing of cohomology in general, it is known that for *any* complex space  $X$  and coherent analytic sheaf  $\mathcal{F}$  on  $X$ ,  $H^p(X, \mathcal{F}) = 0$  for  $p > \dim X$  (cf. [6]), and  $H^p(X, \mathcal{F}) = 0$  for  $p = \dim X$ , at least if  $X$  has no compact  $n$ -dimensional branch [9].

Concerning the converse of Theorem B, it is easy to see that if  $H^1(X, \mathcal{I}) = 0$  for every sheaf of ideals  $\mathcal{I}$  of  $\mathcal{O}$ , where  $\mathcal{O}$  denotes the sheaf of germs of holomorphic functions on  $X$ , then  $X$  is Stein (cf. [4], p. 246). The problem treated here is to decide to what extent the vanishing of the Dolbeault cohomology groups  $H^p(X, \mathcal{O})$ ,  $p > 0$ , guarantees that the space is Stein.

To begin with  $H^p(\mathbb{P}_n, \mathcal{O}) = 0$  for  $p > 0$ , where  $\mathbb{P}_n$  is the  $n$ -dimensional projective space. This is shown in [3], (p. 118), using Hodge theory. Since  $\mathbb{P}_n$  is compact the only global holomorphic functions are the constants and hence they do not separate the points on  $\mathbb{P}_n$  which consequently cannot be Stein. On the other hand it was shown by Laufer [5] that if  $X$  is a Riemann domain over a Stein manifold, such that holomorphic functions separate the points on  $X$ , then  $X$  is Stein if the higher Dolbeault groups vanish. This result was generalized by Siu in [8], who removed the condition of holomorphic separability. He also showed that  $\dim_{\mathbb{C}} H^p(X, \mathcal{O})$ ,  $p > 0$ , cannot even be countable unless  $X$  is Stein.

---

1980 *Mathematics Subject Classification*. Primary 32E10; Secondary 32D05.

*Key words and phrases*. Stein space, domain of holomorphy, Dolbeault cohomology.

As for spaces, it has recently been shown by Fornaess and Narasimhan, [2], that if  $X$  is locally Stein<sup>1</sup>, and relatively compact in a Stein space, then the vanishing of  $H^p(X, \mathcal{O})$  for  $p > 0$  implies that  $X$  is Stein.

In this note we will show, using results from [1], that

(i) if we remove the condition of local Steinness we can still deduce that  $X$  is a domain of holomorphy,

(ii) if  $X$  is relatively compact in a Stein manifold we get under the same conditions that  $X$  is Stein, i.e. we get an alternative proof of (a special case of) Laufer's result.

### Main results

In the following we assume that  $X$  is an open, relatively compact subspace of a Stein space  $S$ . This implies that  $X$  has finite dimension which we denote by  $p$ . Consequently, using a result mentioned above, we have that  $H^q(X, \mathcal{F}) = 0$  for  $q > p$ , for every coherent analytic sheaf  $\mathcal{F}$  on  $X$ .

It may be remarked here that the vanishing of the groups  $H^q(X, \mathcal{F})$  for large  $q$  easily follows using the definition of the Čech cohomology and the fact that  $X$  has finite covering dimension, i.e. every open cover can be refined so that the intersection of any  $2p+2$  open sets is empty.

The reasoning to obtain the following result is more or less standard; we include it for completeness.

A. If  $H^q(X, \mathcal{O}) = 0$  for  $q > 0$ , it follows that for every coherent sheaf  $\mathcal{F}$  on  $S$  we have that  $H^q(X, \mathcal{F}) = 0$  for  $q > 0$ .

PROOF. To begin with we note that  $H^q(X, \mathcal{O}^r) = 0$  for  $q > 0, r > 0$ . This follows immediately using the long exact cohomology sequence corresponding to the short exact sequence

$$0 \rightarrow \mathcal{O} \rightarrow \mathcal{O}^s \rightarrow \mathcal{O}^{s-1} \rightarrow 0$$

and induction.

Let us now take a compact subspace  $K$  of  $S$  such that  $X \subset K$ . Since  $S$  is Stein the holomorphically convex hull  $\hat{K}$  of  $K$  is compact. From Theorem A we then have an exact sequence

$$\mathcal{O}^r_1 \xrightarrow{\varphi_1} \mathcal{F} \rightarrow 0$$

on  $\hat{K}$ . Let  $\mathcal{K}_1$  be the kernel of  $\varphi_1$ . This is a coherent analytic sheaf and the sequence

$$0 \rightarrow \mathcal{K}_1 \rightarrow \mathcal{O}^r_1 \rightarrow \mathcal{F} \rightarrow 0$$

is exact. Using the corresponding long exact sequence we deduce that  $H^q(X, \mathcal{F}) = H^{q+1}(X, \mathcal{K}_1)$  for  $q > 0$ . Repeating the argument with  $\mathcal{K}_1$  replacing  $\mathcal{F}$  and  $\mathcal{K}_2$  denoting the corresponding kernel we obtain  $H^q(X, \mathcal{F}) = H^{q+2}(X, \mathcal{K}_2)$  and so, after  $n$  steps,  $H^q(X, \mathcal{F}) = H^{q+n}(X, \mathcal{K}_n)$ . But the right-hand side vanishes for  $n$  sufficiently large as we saw above, and  $A$  is demonstrated.

Let  $Y$  be a hypersurface in  $S$  defined as the zero-set of a non-constant holomorphic function on  $S$ .

<sup>1</sup> I.e. if every point on the boundary of  $X$  has a neighbourhood  $U$  such that  $U \cap X$  is Stein.

B. Every holomorphic function on  $X \cap Y$  extends to  $X$ .

Since the ideal sheaf  $\mathcal{I}$  defined by  $Y$  is a coherent analytic sheaf on  $S$ , it follows from A that  $H^1(X, \mathcal{I}) = 0$ . The extendability then follows in the usual manner (cf. [4], p. 245).

THEOREM 1. *If  $X$  is an open, relatively compact subspace of the Stein space  $S$ , and  $H^q(X, \mathcal{O}) = 0$  for  $q > 0$ , then  $X$  is a domain of holomorphy.*

PROOF. Let the dimension of  $X$  be  $p$ . If  $p = 2$  the result follows from Theorem 2 of [1] and the fact that all 1-dimensional complex spaces are Stein ([4], p. 270), and hence domains of holomorphy, if we appeal to B above.

Assume that the theorem is proved for  $p < n$ . If  $Y$  is a hypersurface in  $S$  as above, then  $Y$  is Stein,  $\dim(Y \cap X) \leq n - 1$ , and from A we have that  $H^q(Y \cap X, \mathcal{O}) = 0$  for  $q > 0$ . From the induction hypothesis it follows that  $Y \cap X$  is a domain of holomorphy, and using B and Theorem 2 of [1] we deduce that  $X$  is a domain of holomorphy.

THEOREM 2. *Let  $X$  be an open, relatively compact submanifold of the Stein manifold  $S$ . If  $H^q(X, \mathcal{O}) = 0$  for  $q > 0$ , then  $X$  is Stein.*

PROOF. Since the general solvability of the Cousin I problem on a manifold  $Z$  follows from the fact that  $H^1(Z, \mathcal{O}) = 0$ , we see that  $X$  is a Cousin I submanifold of  $S$ , and from A it follows that the same is true for all closed submanifolds of  $X$ , defined as the zero-sets of non-constant holomorphic functions on  $S$ . The result is now a consequence of Theorem 6 in [1].

#### REFERENCES

- [1] BERG, G., Complex spaces with plenty of Stein subvarieties, *Math. Scand.* **51** (1982), 158—162.
- [2] FORNAESS, J. E. and NARASIMHAN, R., The Levi problem on complex spaces with singularities, *Math. Ann.* **248** (1980), 47—72. *MR* **81f**: 32020.
- [3] GRIFFITHS, P. and HARRIS, J., *Principles of algebraic geometry*, J. Wiley, New York, 1978. *MR* **80b**: 14001.
- [4] GUNNING, R. and ROSSI, H., *Analytic functions of several complex variables*, Prentice-Hall, Englewood Cliffs, N. J., 1965. *MR* **31**: #4927.
- [5] LAUFER, H. B., On sheaf cohomology and envelopes of holomorphy, *Ann. of Math.* **84** (1966), 102—118. *MR* **35**: #417.
- [6] REIFFEN, H. J., Riemannsche Hebbarkeitssätze für Cohomologieklassen mit kompaktem Träger, *Math. Ann.* **164** (1966), 272—279. *MR* **33**: #5942.
- [7] SIU, Y. T., A note on Cartan's Theorems A and B, *Proc. Amer. Math. Soc.* **18** (1967), 955—956. *MR* **35**: #6867.
- [8] SIU, Y. T., Non-countable dimensions of cohomology groups of analytic sheaves and domains of holomorphy, *Math. Z.* **102** (1967), 17—29. *MR* **36**: #5394.
- [9] SIU, Y. T., Analytic sheaf cohomology groups of dimension  $n$  of  $n$ -dimensional complex spaces, *Trans. Amer. Math. Soc.* **143** (1969), 77—94. *MR* **40**: #5942.

(Received May 11, 1983)

UPPSALA UNIVERSITET  
MATEMATISKA INSTITUTIONEN  
THUNBERGSVÄGEN 3  
S-752 38 UPPSALA  
SWEDEN





# BERRY—ESSÉEN THEOREMS FOR SIGNED LINEAR RANK STATISTICS UNDER NEAR LOCATION ALTERNATIVES

MUNSUP SEOH and MADAN L. PURI

## Summary

Berry—Esséen theorems for signed linear rank statistics with bounded score generating functions are established under near location alternatives.

## 1. Introduction

Let  $X_{N1}, X_{N2}, \dots, X_{NN}$  be independent r.v.'s (random variables) such that  $X_{Nj} = \Delta\theta_{Nj} + Y_j$ ,  $1 \leq j \leq N$ , where  $\theta_{Nj}$ 's are unknown location parameters,  $\Delta$  is a real parameter,  $\{Y_j\}_{j=1}^\infty$  is a sequence of i.i.d. (independently and identically distributed) r.v.'s with a common c.d.f. (cumulative distribution function)  $F(x)$  and a p.d.f. (probability density function)  $f(x)$ . Let  $R_{Nj}^+$ ,  $1 \leq j \leq N$ , be the rank of  $|X_{Nj}|$  among  $\{|X_{Nk}|: 1 \leq k \leq N\}$ . We consider the signed linear rank statistic

$$(1.1) \quad T_N^+ = \sum_{j=1}^N c_{Nj} a_{NR_{Nj}^+} \operatorname{sgn} X_{Nj}$$

where  $c_{N1}, c_{N2}, \dots, c_{NN}$  are arbitrary regression constants;  $a_{N1}, a_{N2}, \dots, a_{NN}$  are scores; and  $\operatorname{sgn} x = 1$  or  $-1$  according as  $x \geq 0$  or  $x < 0$ .

We assume that the scores  $a_{Nj}$ ,  $1 \leq j \leq N$ , are generated by some known function  $J(t)$  (called a score generating function) defined on the open interval  $(0, 1)$  in either of the following two ways:

$$(1.2) \quad a_{Nj} = EJ(U_{N:j}), \quad j = 1, 2, \dots, N \quad (\text{exact scores})$$

$$(1.3) \quad a_{Nj} = J(EU_{N:j}), \quad j = 1, 2, \dots, N \quad (\text{approximate scores})$$

where  $U_{N:j}$  denotes the  $j$ -th order statistic in the random sample of size  $N$  from the uniform distribution on  $(0, 1)$ .

The asymptotic normality of  $T_N^+$  was established under very general conditions (Hájek [6] and Hušková [9]). It was also shown by Hájek [6] that, when  $f(x)$  is symmetric, the test of the hypothesis  $H: \Delta = 0$  against  $\Delta > 0$  is asymptotically most powerful under suitable assumptions. Recently, to get more precise information than asymptotic normality can provide, the rate of convergence to asymptotic normality and Edgeworth expansions (see e.g. Bickel [2]) has been investigated by several authors. Under the hypothesis  $H$ , when  $f(x)$  is symmetric, Puri and Seoh [15], [16] have derived Berry—Esséen bounds of order  $O(N^{-1/2})$

and Edgeworth expansions with remainder  $o(N^{-1})$  for a wide class of score generating functions including normal scores case. For the corresponding study of similar problems in the case of unsigned linear rank statistic of the form  $T_N = \sum_{j=1}^N c_{Nj} a_{NR_{Nj}}$ , where  $R_{Nj}$  is the rank of  $X_{Nj}$  among  $\{X_{Nk}; 1 \leq k \leq N\}$ , the reader is referred to Hájek [6], [7] and Dupáček and Hájek [4] (for asymptotic normality); Jurečková and Puri [12], Bergström and Puri [1], Hušková [10], [11] (for the rate of convergence); and Does [3] (for the rate of convergence and Edgeworth expansions), among others.

Recently, Puri and Wu [17] have derived a bound of order  $O(N^{-1/2+\delta})$ ,  $\delta > 0$ , under the null hypothesis, viz.  $\Delta = 0$ , and under contiguous location alternatives, for the rate of convergence to normality of the statistic (1.1) with approximate scores assuming that the score generating function satisfies Lipschitz's condition of order one.

In this paper, we derive the Berry—Essén bound  $O(N^{-1/2})$  for the statistic (1.1) with exact as well as approximate scores for the case of near location alternatives satisfying Assumption B below. Our assumptions on  $c$ 's,  $\theta$ 's and  $F$  are relatively milder than those in Puri and Wu [17]. On the other hand, our assumption on the score generating function  $J$  is slightly stronger than that in Puri and Wu (op. cit.).

## 2. Assumptions and main theorems

Throughout this paper, we make the following assumptions.

**ASSUMPTION A.** The regression constants  $c_{N1}, c_{N2}, \dots, c_{NN}$  satisfy

$$\sum_{j=1}^N c_{Nj}^2 = 1, \quad \sum_{j=1}^N |c_{Nj}|^3 = O(N^{-1/2}).$$

**ASSUMPTION B.** The location parameters  $\theta_{N1}, \theta_{N2}, \dots, \theta_{NN}$  satisfy

$$\sum_{j=1}^N |\theta_{Nj}|^3 = O(N^{-1/2}).$$

**ASSUMPTION C.** The p.d.f.  $f(x)$  is symmetric about zero and has a Radon—Nikodym derivative  $f'(x)$  with respect to the Lebesgue measure such that

$$\int_{-\infty}^{\infty} |f'(x)| dx < \infty.$$

**ASSUMPTION D.** The score generating function  $J$  is differentiable on  $(0, 1)$ ,  $\int_0^1 J^2(t) dt = 1$  and the derivative  $J'$  satisfies Lipschitz's condition of order  $\delta$ ,  $0 < \delta \leq 1$ ; i.e.

$$|J'(s) - J'(t)| \leq \Delta |s - t|^\delta, \quad 0 < s, t < 1,$$

for some positive constant  $\Delta$ .

Defining for each  $N$ ,  $N \geq 1$ ,

$$(2.1) \quad \tau_N^2 = \sigma^2(T_N^+) = \text{Var } T_N^+, \quad T_N^* = \tau_N^{-1}(T_N^+ - ET_N^+),$$

we state our main theorems.

**THEOREM 2.1.** *If Assumptions A through D are satisfied, then as  $N \rightarrow \infty$*

$$\sup_x |\mathbb{P}(T_N^* \leq x) - \Phi(x)| = O(N^{-1/2})$$

where  $\Phi(x)$  is the standard normal c.d.f.

We need a stronger assumption to get simple normalizing constants by removing terms of order  $O(N^{-1/2})$  in asymptotic expansions of  $ET_N^+$  and  $\text{Var } T_N^+$ .

**ASSUMPTION E.** The second Radon—Nikodym derivative  $f''(x)$  exists and

$$\int_{-\infty}^{\infty} |f''(x)| dx < \infty.$$

**THEOREM 2.2.** *Suppose Assumptions A through E are satisfied. Then as  $N \rightarrow \infty$ ,*

$$(2.2) \quad \sup_x |\mathbb{P}(T_N^+ - \mu_N \leq x) - \Phi(x)| = O(N^{-1/2})$$

where

$$\mu_N = \sum_{j=1}^N c_{Nj} \int_{-\infty}^{\infty} J(F^+(|x|)) \operatorname{sgn} x (f(x - \theta_{Nj}) - f(x)) dx$$

and  $F^+(x)$  is the c.d.f. of the r.v.  $|Y_1|$ .

**REMARK 2.1.** If the third integrable Radon—Nikodym derivative exists, then

$$\begin{aligned} \mu_N &= \sum_{j=1}^N c_{Nj} \theta_{Nj} \int_{-\infty}^{\infty} J(F^+(|x|)) \operatorname{sgn} x f'(x) dx + O(N^{-1/2}) = \\ &= \sum_{j=1}^N c_{Nj} \theta_{Nj} \int_0^1 J(t) \frac{f'((F^+)^{-1}(t))}{f((F^+)^{-1}(t))} dt + O(N^{-1/2}). \end{aligned}$$

**REMARK 2.2.** Assumptions A and B are well-known and quite satisfactory for practical purposes. Note that Assumption B implies that  $\sum_{j=1}^N \theta_{Nj}^2 = O(1)$ .

**REMARK 2.3.** Assumption C implies that  $f$  is absolutely continuous and bounded.

Furthermore  $\lim_{x \rightarrow \pm \infty} f(x) = 0$  and  $\int_{-\infty}^{\infty} f'(x) dx = 0$  (see e.g. Hájek and Šidák [8], p. 20).

**REMARK 2.4.** Assumption D is rather restrictive. It requires that the score generating functions be bounded. Thus it does not cover the normal scores case. But note that our Assumption D is weaker than one of assuming bounded second derivatives.

We note that we can replace  $\Delta\theta_{Nj}$  by simply  $\theta_{Nj}$  without loss of generality. From now on, we shall omit indices  $N$  in  $c_{Nj}$ ,  $a_{Nj}$ ,  $R_{Nj}$ ,  $X_{Nj}$ ,  $\theta_{Nj}$ , etc., whenever it causes no confusion. The constant  $\delta$  is reserved to denote the order of Lipschitz's condition (c.f. Assumption D). We shall also omit bounds of integration  $\int_{-\infty}^{\infty} f(x)dx$  etc. if integration is taken over the whole real line.

### 3. Three lemmas

In this section we prove three lemmas which play important roles in the proofs of our main theorems.

**LEMMA 3.1.** *Suppose Assumption C is satisfied. Let  $v(x_1, x_2, \dots, x_N)$  be a function of  $N$  independent variables  $x_1, x_2, \dots, x_N$  and define a r.v.  $V_N = v(|X_1|, |X_2|, \dots, |X_N|)$ . If  $E|V_N| < \infty$ , then we have, for distinct indices  $j_1, j_2, \dots, j_n$ ,  $1 \leq n \leq N$ ,*

$$(3.1) \quad \begin{aligned} E(V_N \prod_{v=1}^n \operatorname{sgn} X_{j_v}) &= E[E\{V_N | X_{j_1}, X_{j_2}, \dots, X_{j_n}\} \prod_{v=1}^n \operatorname{sgn} X_{j_v}] = \\ &= \int \dots \int E\{V_N | X_{j_1} = y_1, \dots, X_{j_n} = y_n\} \prod_{v=1}^n \{\operatorname{sgn} y_v (f(y_v - \theta_{j_v}) - f(y_v))\} dy_v. \end{aligned}$$

Furthermore, if there is constant  $K_N$  (depending only on  $N$ ) such that  $|E\{V_N | X_{j_1}, X_{j_2}, \dots, X_{j_n}\}| \leq K_N$ , we have

$$(3.2) \quad |E(V_N \prod_{v=1}^n \operatorname{sgn} X_{j_v})| \leq K_N \prod_{v=1}^n |\theta_{j_v}| \left( \int |f'(t)| dt \right)^n.$$

**PROOF.** Noting that  $E\{V_N | X_{j_1}, \dots, X_{j_n}\}$  is a function of  $|X_{j_1}|, \dots, |X_{j_n}|$ , we set

$$(3.3) \quad \tilde{v}(|y_1|, |y_2|, \dots, |y_n|) = E\{V_N | X_{j_1} = y_1, \dots, X_{j_n} = y_n\}.$$

Since the vector  $(\operatorname{sgn} Y_1, \operatorname{sgn} Y_2, \dots, \operatorname{sgn} Y_n)$ , of which components are i.i.d. with zero means, is independent of the vector  $(|Y_1|, |Y_2|, \dots, |Y_n|)$  (see e.g. Hájek and Šidák [8]), we have

$$(3.4) \quad \begin{aligned} \int \dots \int \tilde{v}(|y_1|, \dots, |y_n|) \prod_{v=1}^n \{\operatorname{sgn} y_v f(y_v) dy_v\} &= \\ &= E\{\tilde{v}(|Y_1|, |Y_2|, \dots, |Y_n|) \prod_{v=1}^n \operatorname{sgn} Y_v\} = 0. \end{aligned}$$

Also it follows by Fubini's theorem that

$$(3.5) \quad \int \dots \int \tilde{v}(|y_1|, \dots, |y_n|) \prod_{v=1}^n \{\operatorname{sgn} y_v g(y_v) dy_v\} = 0$$

if  $g(y_v) = f(y_v - \theta_{j_v})$  or  $= f(y_v)$ ,  $1 \leq v \leq n$ , and the latter equality holds at least once. Hence (3.1) follows by (3.3), (3.4) and (3.5).

We now turn to (3.2). Applying Fubini's theorem again, we obtain

$$(3.6) \quad \int |f(x-\theta) - f(x)| dx \leq \int_{-\infty}^{\infty} \int_x^{x-\theta} |f'(t)| dt dx \leq |\theta| \int |f'(t)| dt$$

which, together with (3.1), implies (3.2). Thus the proof is complete.

LEMMA 3.2. *Under assumptions of Theorem 2.1, we have*

$$E \left( \sum_{j=1}^N c_j [\tilde{a}_{R_j^+} - J(R_j^+/(N+1))] \operatorname{sgn} X_j \right)^2 = O(N^{-1-\delta})$$

where

$$\tilde{a}_j = EJ(U_{N,j}), \quad 1 \leq j \leq N.$$

PROOF. Because of Assumption D, Taylor's expansion yields

$$\tilde{a}_j = J\left(\frac{j}{N+1}\right) + E\left(U_{N,j} - \frac{j}{N+1}\right) J'\left(\lambda_j U_{N,j} + (1-\lambda_j) \frac{j}{N+1}\right)$$

where  $0 \leq \lambda_j \leq 1$ ,  $1 \leq j \leq N$ . Since  $E(U_{N,j} - j/(N+1))J'(j/(N+1)) = 0$  and  $J'$  satisfies Lipschitz's condition of order  $\delta$ , we derive uniformly for  $1 \leq j \leq N$

$$(3.7) \quad \tilde{a}_j - J\left(\frac{j}{N+1}\right) = O\left(E\left|U_{N,j} - \frac{j}{N+1}\right|^{1+\delta}\right) = O(N^{-1/2-\delta/2})$$

with the aid of Hölder's inequality. Applying Lemma 3.1, we obtain

$$(3.8) \quad \begin{aligned} E \left( \sum_{j=1}^N c_j \left\{ \tilde{a}_{R_j^+} - J\left(\frac{R_j^+}{N+1}\right) \right\} \operatorname{sgn} x_j \right)^2 &= \sum_{j=1}^N c_j^2 E \left\{ \tilde{a}_{R_j^+} - J\left(\frac{R_j^+}{N+1}\right) \right\}^2 + \\ &+ \sum_{j \neq k} \sum c_j c_k E \left[ \left\{ \tilde{a}_{R_j^+} - J\left(\frac{R_j^+}{N+1}\right) \right\} \left\{ \tilde{a}_{R_k^+} - J\left(\frac{R_k^+}{N+1}\right) \right\} \operatorname{sgn} X_j \operatorname{sgn} X_k \right] = \\ &= \sum_{j=1}^N c_j^2 O(N^{-1-\delta}) + \sum_{j \neq k} \sum |c_j c_k \theta_j \theta_k| O(N^{-1-\delta}). \end{aligned}$$

Finally, Cauchy—Schwartz's inequality, Assumptions A and B complete the proof.

We define r.v.'s for  $1 \leq j \leq N$

$$(3.9) \quad \varrho_j = R_j^+/(N+1) = (N+1)^{-1} \left\{ 1 + \sum_{k \neq j} u(|X_j| - |X_k|) \right\},$$

$$(3.10) \quad \varrho_{jj} = E(\varrho_j | X_j) = (N+1)^{-1} \left\{ 1 + \sum_{k \neq j} G_k^+(|X_j|) \right\},$$

where  $u(x) = (\operatorname{sgn} x + 1)/2$  and  $G_k^+(x)$  is the c.d.f. of the r.v.  $|X_k|$ , i.e.,

$$(3.11) \quad G_k^+(x) = F(x - \theta_k) - F(-x - \theta_k).$$

LEMMA 3.3. Let  $v$  and  $r$  be positive integers such that  $v \leq N-1$ . Then for  $1 \leq j \leq N$  we have that

$$(3.12) \quad E\{(\varrho_j - \varrho_{jj})^{2r} | X_j, X_{l_1}, \dots, X_{l_v}\} \leq N^{-2r} \{(2v)^{2r} + (16er(N-1-v))^r\}$$

where  $j, l_1, \dots, l_v$  are distinct indices among  $\{1, 2, \dots, N\}$ .

PROOF. Denote

$$(3.13) \quad h(X_j, X_k) = u(|X_j| - |X_k|) - G_k^+(|X_j|), \quad 1 \leq j, k \leq N,$$

and then, in view of (3.9) and (3.10), we can write

$$(3.14) \quad \varrho_j - \varrho_{jj} = (N+1)^{-1} \sum_{k \neq j}^N h(X_j, X_k) = (N+1)^{-1} (Z_1 + Z_2)$$

where the r.v.'s  $Z_1$  and  $Z_2$  are given by

$$(3.15) \quad Z_1 = \sum_{k=1}^v h(X_j, X_{l_k}), \quad Z_2 = \sum_{\substack{k \neq j \\ k \neq l_1, \dots, l_v}}^N h(X_j, X_k).$$

Then it follows by Hölder's inequality that

$$(3.16) \quad \begin{aligned} & E\{(\varrho_j - \varrho_{jj})^{2r} | X_j, X_{l_1}, X_{l_2}, \dots, X_{l_v}\} \leq \\ & \leq N^{-2r} 2^{2r} [E\{Z_1^{2r} | X_j, X_{l_1}, \dots, X_{l_v}\} + E\{Z_2^{2r} | X_j, X_{l_1}, \dots, X_{l_v}\}]. \end{aligned}$$

Conditionally given  $X_j, X_{l_1}, \dots, X_{l_v}$ ,  $Z_2$  is the sum of independent r.v.'s with zero means and thus we may apply Lemma 3.1 of Seoh, Ralescu and Puri [20] to obtain that

$$(3.17) \quad E\{Z_2^{2r} | X_j, X_{l_1}, \dots, X_{l_v}\} \leq (4er)^r (N-1-v)^r.$$

Since  $|Z_1| \leq v$ , (3.16) and (3.17) ensure (3.12) to complete the proof.

#### 4. Proofs of theorems

A standard approach to establish Berry—Esséen theorems is to invoke Esséen's smoothing lemma (see e.g. Feller ([5], p. 538)), which implies that for all  $\gamma > 0$

$$(4.1) \quad \sup |P(T_N^* \leq x) - \Phi(x)| \leq \frac{1}{\pi} \int_{|t| \leq \gamma N^{1/2}} |t|^{-1} |\psi_N^*(t) - e^{-t^2/2}| dt + O(N^{-1/2})$$

where  $\psi_N^*(t)$  is the characteristic function of  $T_N^*$ .

To estimate the integral in (4.1), being the proper order  $O(N^{-1/2})$ , we shall derive an asymptotic expression of  $\psi_N^*(t)$  on the range  $|t| \leq \log N$  and find a sharp bound on the integrand on the remaining range  $\log N \leq |t| \leq \gamma N^{1/2}$ .

PROOF OF THEOREM 2.1. First we consider the approximate scores. Recall that the r.v.'s  $\varrho_j$  and  $\varrho_{jj}$ , and the c.d.f.  $G_j^+(x)$  of the r.v.  $|X_j|$ ,  $1 \leq j \leq N$ , are defined by (3.9), (3.10) and (3.11), respectively.



In view of Assumption C, Taylor's expansion yields that for any  $x$  and  $\theta$

$$(4.2) \quad F(x-\theta) = F(x) - \theta f(x) + \int_x^{x-\theta} (x-\theta-t)f'(t) dt.$$

Since  $f(x) = f(-x)$  for all  $x$ , we obtain from (3.11) and (4.2) that for  $1 \leq k \leq N$

$$(4.3) \quad G_k^+(x) = F^+(x) + R(x, \theta_k)$$

where  $F^+(x) = F(x) - F(-x)$  is the c.d.f. of  $|Y_1|$  and

$$(4.4) \quad R(x, \theta) = \int_x^{x-\theta} (x-\theta-t)f'(t) dt + \int_{-x}^{-x-\theta} (x+\theta+t)f'(t) dt.$$

Furthermore it follows from (3.10) and (4.3) that uniformly for  $1 \leq j \leq N$

$$(4.5) \quad \varrho_{jj} = F^+(|X_j|) + (N+1)^{-1} \sum_{k \neq j}^N R(|X_j|, \theta_k) + O(N^{-1}).$$

Let  $S_N^+$  be the first two terms of Taylor's expansion of  $T_N^+$ , viz.

$$(4.6) \quad S_N^+ = \sum_{j=1}^N c_j \{J(\varrho_{jj}) + (\varrho_j - \varrho_{jj})J(\varrho_{jj})\} \operatorname{sgn} X_j$$

and define

$$(4.7) \quad I_1 = \sum_{j=1}^N c_j \{J(\varrho_{jj}) \operatorname{sgn} X_j - E(J(\varrho_{jj}) \operatorname{sgn} X_j),$$

$$I_2 = \frac{1}{N+1} \sum_{j \neq k}^N c_j J'(\varrho_{jj}) \{u(|X_j| - |X_k|) - G_k^+(|X_j|)\} \operatorname{sgn} X_j,$$

so that  $S_N^+ - ES_N^+ = I_1 + I_2$ . To approximate  $T_N^+$  by  $S_N^+$ , we need

LEMMA 4.1. *Under the assumptions of Theorem 2.1, we have*

$$(4.8) \quad \sigma^2(T_N^+ - S_N^+) = O(N^{-1-\delta}).$$

PROOF. Taylor's expansion yields that with real numbers  $0 \leq \lambda_j \leq 1$ ,  $1 \leq j \leq N$ ,

$$T_N^+ = \sum_{j=1}^N c_j \{J(\varrho_{jj}) + (\varrho_j - \varrho_{jj})J'(\lambda_j \varrho_j + (1-\lambda_j)\varrho_{jj})\} \operatorname{sgn} X_j.$$

Setting  $W_j = (\varrho_j - \varrho_{jj})\{J'(\lambda_j \varrho_j + (1-\lambda_j)\varrho_{jj}) - J'(\varrho_{jj})\}$ ,  $1 \leq j \leq N$ , we have

$$(4.9) \quad E(T_N^+ - S_N^+)^2 = \sum_{j=1}^N c_j^2 E W_j^2 + \sum_{j \neq k}^N c_j c_k E(W_j W_k \operatorname{sgn} X_j \operatorname{sgn} X_k).$$

Applying Lemma 3.3, we obtain that for any positive integer  $r$

$$(4.10) \quad E\{|\varrho_j - \varrho_{jj}|^{2r} |X_j, X_k\} \leq CN^{-r}$$

where  $C$  is a constant depending only on  $r$ . Then Hölder's inequality ensures (4.10) for any real  $r$ . It follows by Lipschitz's condition on  $J'$  and (4.10) that

$$(4.11) \quad \mathbb{E}W_j^2 = O(\mathbb{E}|\varrho_j - \varrho_{jj}|^{2(1+\delta)}) = O(N^{-1-\delta}).$$

Thus, in view of Assumption A, the first sum on the R.H.S. (right-hand side) of (4.9) is of order  $O(N^{-1-\delta})$ .

We next estimate the second sum on the R.H.S. of (4.9). It follows from (4.11) and Hölder's inequality that uniformly for  $1 \leq j, k \leq N$

$$(4.12) \quad \mathbb{E}\{W_j W_k | X_j, X_k\} \leq \{\mathbb{E}(W_j^2 | X_j, X_k) \mathbb{E}(W_k^2 | X_j, X_k)\}^{1/2} = O(N^{-1-\delta}).$$

Hence, an application of Lemma 3.1, Assumptions A and B, and (4.12) ensure that

$$(4.13) \quad \sum_{j \neq k} c_j c_k \mathbb{E}W_j W_k \operatorname{sgn} X_j \operatorname{sgn} X_k = O(N^{-1-\delta} \sum_{j \neq k} \sum |c_j c_k \theta_j \theta_k|) = O(N^{-1-\delta}).$$

Thus the proof is complete.

We now compute a number of moments of r.v.'s  $I_1$  and  $I_2$ , defined by (4.7).

LEMMA 4.2. *Under assumptions of Theorem 2.1, we have*

$$(4.14) \quad \mathbb{E}I_1^2 = 1 + \sum_{j=1}^N c_j^2 \int J^2(F^+(|x|))(f(x - \theta_j) - f(x)) dx + O(N^{-1/2-1/6}),$$

$$(4.15) \quad \mathbb{E}I_2^2 = O(N^{-1}), \quad \mathbb{E}I_1 I_2 = O(N^{-1}).$$

PROOF. Because of (4.5) and Assumption D, we obtain

$$(4.16) \quad J(\varrho_{jj}) = J(F^+(|X_j|)) + O(N^{-1} \sum_{k=1}^N |R(|X_j|, \theta_k)|) + O(N^{-1})$$

uniformly in  $1 \leq j \leq N$ . Set  $\|f\| = \sup_x |f(x)|$  and let  $\tilde{f}_j(x)$  be the p.d.f. of  $|X_j|$ . Then, for any  $x \geq 0$ ,  $|\tilde{f}_j(x)| = |f(x - \theta_j) + f(-x - \theta_j)| \leq 2\|f\| < \infty$ . Put  $K = 2 \int |f'(t)| dt$  and then it is clear from (4.4) that

$$(4.17) \quad |R(x, \theta)| \leq \theta K$$

which, together with Assumption B, ensures that

$$(4.18) \quad (N^{-1} \sum_{k=1}^N |R(x, \theta_k)|)^2 \leq K^2 N^{-1}.$$

Furthermore, by elementary computations based on Fubini's theorem we get

$$(4.19) \quad \mathbb{E}|R(|X_j|, \theta)| \leq 2K\|f\|\theta^2, \quad 1 \leq j \leq N.$$

Thus it follows from (4.16), (4.18), (4.19) and Lemma 3.1 that uniformly for  $1 \leq j \leq N$

$$(4.20) \quad \begin{aligned} \mathbb{E}J^2(\varrho_{jj}) &= \mathbb{E}J^2(F^+(|X_j|)) + O(N^{-1}), \\ \mathbb{E}J(\varrho_{jj}) \operatorname{sgn} X_j &= \mathbb{E}J(F^+(|X_j|)) \operatorname{sgn} X_j + O(N^{-1}) = O(|\theta_j|) + O(N^{-1}). \end{aligned}$$

As  $I_1$  is a sum of independent r.v.'s with zero means, we derive

$$\begin{aligned}
 EI_1^2 &= \sum_{j=1}^N c_j^2 \{EJ^2(\varrho_{jj}) - (EJ(\varrho_{jj}) \operatorname{sgn} X_j)^2\} = \\
 (4.21) \quad &= \sum_{j=1}^N c_j^2 \left\{ \int J^2(F^+(|x|)) f(x - \theta_j) dx + O(\theta_j^2) + O(N^{-1}) \right\} = \\
 &= 1 + \sum_{j=1}^N c_j^2 \int J^2(F^+(|x|)) (f(x - \theta_j) - f(x)) dx + O(N^{-1/2-1/6})
 \end{aligned}$$

because of (4.20), Assumptions A and B, and  $\int_0^1 J^2(t) dt = 1$ .

We now turn to (4.15). Since  $I_2 = \sum_{j=1}^N c_j(\varrho_j - \varrho_{jj})J'(\varrho_{jj}) \operatorname{sgn} X_j$ , applying (4.10) and Lemma 3.1, we derive

$$\begin{aligned}
 EI_2^2 &= \sum_{j=1}^N c_j^2 E((\varrho_j - \varrho_{jj})J'(\varrho_{jj}))^2 + \\
 (4.22) \quad &+ \sum_{j \neq k} c_j c_k E(\varrho_j - \varrho_{jj})(\varrho_k - \varrho_{kk})J'(\varrho_{jj})J'(\varrho_{kk}) \operatorname{sgn} X_j \operatorname{sgn} X_k = \\
 &= O(N^{-1}) + \sum_{j \neq k} |c_j c_k \theta_j \theta_k| O(E\{| \varrho_j - \varrho_{jj} | | \varrho_k - \varrho_{kk} | | X_j, X_k \}) = O(N^{-1}).
 \end{aligned}$$

Similarly we obtain by (easy) direct computations that  $E I_1 I_2 = O(N^{-1})$  and the proof is complete.

It follows from (3.6) and Lemma 4.2 that

$$(4.23) \quad \sigma^2(I_1) = EI_1^2 = 1 + O(N^{-1/6}),$$

$$(4.24) \quad \sigma_N^2 = \operatorname{Var} S_N^+ = \sigma^2(I_1) + O(N^{-1}),$$

which, together with Lemma 4.1, ensure that

$$\begin{aligned}
 (4.25) \quad \tau_N^2 &= \operatorname{Var} T_N^+ = \sigma_N^2 + 2 \operatorname{Cov}(S_N^+, T_N^+ - S_N^+) + \sigma^2(T_N^+ - S_N^+) = \\
 &= \sigma_N^2 + O(N^{-1/2-\delta/2}) = 1 + O(N^{-1/6}).
 \end{aligned}$$

We next consider the behaviour of the characteristic function  $\psi_N^*$  for large values of the argument.

**LEMMA 4.3.** *Suppose that assumptions of Theorem 2.1 are satisfied. Then there exist positive numbers  $B$ ,  $\beta$  and  $\gamma$  such that*

$$(4.26) \quad |\psi_N^*(t)| = |E e^{itT_N^*}| \leq B N^{-\beta \log N}$$

for  $\log N \leq |t| \leq \gamma N^{1/2}$ ,  $N \geq 1$ .

**PROOF.** Since this lemma is a special case of Corollary 2.2 of Seoh [19], we need to check that assumptions of the corollary are satisfied. Define

$$(4.27) \quad \hat{T}_N = \sum_{j=1}^N \hat{c}_j a_{R_j^+} \operatorname{sgn} X_j$$

where  $\hat{c}_j = N^{1/2} c_j \tau_N^{-1}$ ,  $1 \leq j \leq N$ . Then we have

$$(4.28) \quad \psi_N^*(t) = E \exp \{it N^{-1/2} (\hat{T}_N - E\hat{T}_N)\}.$$

Because of Assumptions A and D and (4.25), it is easy to check (cf. Lemma 3.2 of Puri and Seoh [15]) that for some positive numbers  $c, C, a$  and  $A$

$$(4.29) \quad \sum_{j=1}^N |\hat{c}_j| \geq cN, \quad \sum_{j=1}^N \hat{c}_j^2 \leq CN,$$

$$\sum_{j=1}^N |a_j| \geq aN, \quad \sum_{j=1}^N a_j^2 \leq AN.$$

We next take  $\tilde{\theta} = \max_{1 \leq j \leq N} |\theta_j|$  and then

$$(4.30) \quad \begin{aligned} & \int \left\{ \frac{1}{N} \sum_{j=1}^N (f(x - \theta_j) - f(x))^2 \right\}^{1/2} dx = \int \left\{ \frac{1}{N} \sum_{j=1}^N \left( \int_x^{x-\theta_j} f'(y) dy \right)^2 \right\}^{1/2} dx \leq \\ & \leq \int_{-\infty}^{\infty} \int_{x-\tilde{\theta}}^{x+\tilde{\theta}} |f'(y)| dy dx = 2\tilde{\theta} \int |f'(y)| dy = O(N^{-1/6}). \end{aligned}$$

It follows from Section 3 of van Zwet [21] that (4.30) implies that

$$(4.31) \quad \sum_{j=1}^N \int \frac{(f(x - \theta_j) - f(x))^2}{f(x)} dx = O(N).$$

Hence (4.28), (4.29), (4.31) and our model of near location alternatives ensure that assumptions of the corollary are satisfied. Thus an application completes the proof.

Because of Lemma 4.3, we have

$$(4.32) \quad \int_{\log N \leq |t| \leq \gamma N^{1/2}} |t|^{-1} |\psi_N^*(t) - e^{-t^2/2}| dt = O(N^{-1/2})$$

which, together with (4.1), ensures that it is sufficient to show

$$(4.33) \quad \int_{|t| \leq \log N} |t|^{-1} |\psi_N^*(t) - e^{-t^2/2}| dt = O(N^{-1/2})$$

in order to prove Theorem 2.1.

Define  $S_N^* = \sigma_N^{-1} (S_N^+ - ES_N^+)$ ,  $N \geq 1$ . Then

$$(4.34) \quad \begin{aligned} E(T_N^* - S_N^*)^2 &= \frac{2\tau_N \sigma_N - 2 \text{Cov}(T_N^+, S_N^+)}{\tau_N \sigma_N} = \frac{\sigma^2(T_N^+ - S_N^+) - (\tau_N - \sigma_N)^2}{\tau_N \sigma_N} \\ &\leq (\tau_N \sigma_N)^{-1} \sigma^2(T_N^+ - S_N^+) = O(N^{-1-\delta}) \end{aligned}$$

in view of (4.8), (4.24) and (4.25). Hence we have

$$(4.35) \quad |\psi_N^*(t) - Ee^{itS_N^*}| = O(|t| E|T_N^* - S_N^*|) = O(|t| N^{-1/2-\delta/2}).$$

Finally we need

LEMMA 4.4. Under assumptions of Theorem 2.1, we have

$$(4.36) \quad \mathbb{E} I_2 e^{it\sigma_N^{-1} I_1} = O(|t| N^{-1/2-1/6} P(t))$$

where  $P(t)$  is a fixed polynomial.

PROOF. Define  $\tilde{J}(X_j) = J(\varrho_{jj}) \operatorname{sgn} X_j - \mathbb{E} J(\varrho_{jj}) \operatorname{sgn} X_j$ . Then we have

$$(4.37) \quad \begin{aligned} |\mathbb{E} I_2 e^{it\sigma_N^{-1} I_1}| &= |(N+1)^{-1} \sum_{j \neq k} c_j \left[ \prod_{l \neq j, k} \mathbb{E} \exp(it\sigma_N^{-1} c_l \tilde{J}(X_l)) \right] \times \\ &\quad \times \mathbb{E} [\exp \{it\sigma_N^{-1} (c_j \tilde{J}(X_j) + c_k \tilde{J}(X_k))\} J'(\varrho_{jj}) h(X_j, X_k) \operatorname{sgn} X_j] \leq \\ &\leq N^{-1} \sum_{j \neq k} |c_j| \left[ |\mathbb{E} J'(\varrho_{jj}) h(X_j, X_k) \operatorname{sgn} X_j \{1 + it\sigma_N^{-1} (c_j \tilde{J}(X_j) + c_k \tilde{J}(X_k)) + \right. \\ &\quad \left. + 1/2(it\sigma_N^{-1})^2 (c_j \tilde{J}(X_j) + c_k \tilde{J}(X_k))^2\}| + O(|t|^3(|c_j|^3 + |c_k|^3)) \right], \end{aligned}$$

where  $h(X_j, X_k) = u(|X_j| - |X_k|) - G_k^+(|X_j|)$ . Note that

$$(4.38) \quad \mathbb{E} J'(\varrho_{jj}) h(X_j, X_k) \operatorname{sgn} X_j \{1 + it\sigma_N^{-1} c_j \tilde{J}(X_j) + (it\sigma_N^{-1} c_j \tilde{J}(X_j))^2\} = 0.$$

Furthermore, applying Lemma 3.1, we derive uniformly in  $j$  and  $k$

$$(4.39) \quad \begin{aligned} &|\mathbb{E} J'(\varrho_{jj}) h(X_j, X_k) \operatorname{sgn} X_j \{c_k \tilde{J}(X_k) + 2c_j c_k \tilde{J}(X_j) \tilde{J}(X_k) + c_k^2 \tilde{J}^2(X_k)\}| = \\ &= |\mathbb{E} J'(\varrho_{jj}) h(X_j, X_k) \operatorname{sgn} X_j \{c_k J(\varrho_{kk}) + 2c_j c_k J(\varrho_{jj}) J(\varrho_{kk}) \operatorname{sgn} X_j + \\ &\quad + 2c_j c_k (\mathbb{E} J(\varrho_{jj}) \operatorname{sgn} X_j) J(\varrho_{kk}) + 2c_k^2 J(\varrho_{kk})\} \operatorname{sgn} X_k + c_k^2 J^2(\varrho_{kk})| = \\ &= O(|c_k \theta_j \theta_k| + |c_j c_k \theta_k| + |c_j c_k \theta_j \theta_k| + |c_k^2 \theta_j \theta_k| + |c_k^2 \theta_j|). \end{aligned}$$

Because of Assumptions A and B, it follows that

$$(4.40) \quad \begin{aligned} \sum_{j \neq k} |c_j c_k \theta_j \theta_k| &= O(1), \quad \sum_{j \neq k} c_j^2 |c_k \theta_k| = O(1), \\ \sum_{j \neq k} |c_j \theta_j| c_k^2 &= O(1), \quad \sum_{j \neq k} c_j^4 = O(N^{1/2-1/6}), \\ \sum_{j \neq k} |c_j c_k^3| &= O(1). \end{aligned}$$

Thus the proof is complete by combining (4.37) through (4.40).

It follows from (4.15), (4.36) and Lemma XV. 4.1 of Feller [5] that

$$(4.41) \quad |\mathbb{E} e^{itS_N^*} - \mathbb{E} e^{it\sigma_N^{-1} I_1}| = O(|t| |\mathbb{E} e^{it\sigma_N^{-1} I_1} I_2|) + O(t^2 \mathbb{E} I_2^2) = O(|t| N^{-1/2-1/6} P(t))$$

where  $P(t)$  is a fixed polynomial.

Next task is clearly to evaluate the leading term  $\mathbb{E} e^{it\sigma_N^{-1} I_1}$ . Note that  $I_1 = \sum_{j=1}^N c_j \{J(\varrho_{jj}) \operatorname{sgn} X_j - \mathbb{E} J(\varrho_{jj}) \operatorname{sgn} X_j\}$  is a sum of independent r.v.'s with zero means

and finite absolute moments of any order such that

$$(4.42) \quad \sum_{j=1}^N \mathbb{E} |c_j \{J(\varrho_{jj}) \operatorname{sgn} X_j - \mathbb{E} J(\varrho_{jj}) \operatorname{sgn} X_j\}|^3 = O(N^{-1/2}).$$

An application of Lemma V.2.1 of Petrov [14] yields the existence of positive number  $\alpha$  such that

$$(4.43) \quad |\mathbb{E} e^{it\sigma^{-1}(I_1)I_1} - e^{-t^2/2}| = O(N^{-1/2} |t|^3 e^{-t^2/3})$$

uniformly for  $|t| \leq \alpha N^{1/2}$ . Replacing  $t$  by  $\sigma(I_1)\sigma_N^{-1}t$  and expanding  $\exp(-\sigma^2(I)\sigma_N^{-2}t^2/2)$ , we obtain that uniformly for  $|t| \leq \log N$

$$(4.44) \quad \begin{aligned} |\mathbb{E} e^{it\sigma_N^{-1}I_1} - e^{-t^2/2}| &= O(t^2 |\sigma_N^2 - \sigma^2(I_1)|) + O(N^{-1/2} |t|^3 e^{-\theta t^2}) = \\ &= O(N^{-1} t^2) + O(N^{-1/2} |t|^3 e^{-\theta t^2}) \end{aligned}$$

where  $0 < \theta < 1/2$ . It follows from (4.35), (4.41) and (4.44) that uniformly for  $|t| \leq \log N$

$$(4.45) \quad |\psi_N^*(t) - e^{-t^2/2}| = O(|t| N^{-1/2-\varepsilon} P(t)) + O(N^{-1/2} |t|^3 e^{-\theta t^2})$$

where  $\varepsilon = \min(\delta/2, 1/6)$ ,  $0 < \theta < 1/2$  and  $P(t)$  is a fixed polynomial. Since (4.45) ensures (4.33), the proof is complete for approximate scores.

We now consider exact scores. To distinguish statistics with exact scores and approximate scores, we introduce additional notations. Denoting exact scores by  $\tilde{a}_j = \mathbb{E} J(U_{N,j})$ ,  $1 \leq j \leq N$ , we put

$$(4.46) \quad \begin{aligned} \tilde{T}_N^+ &= \sum_{j=1}^N c_j \tilde{a}_{R_j^+} \operatorname{sgn} X_j, \quad \tilde{\tau}_N^2 = \operatorname{Var} \tilde{T}_N^+ \\ \tilde{T}_N^* &= \tilde{\tau}_N^{-1} (\tilde{T}_N^+ - \mathbb{E} \tilde{T}_N^+), \quad \tilde{\psi}_N^*(t) = \mathbb{E} e^{it\tilde{T}_N^*}. \end{aligned}$$

For the statistic  $T_N^+ = \sum_{j=1}^N c_j J(R_j^+/(N+1)) \operatorname{sgn} X_j$ , we use same notations as before.

It follows from Lemma 3.2 and (4.25) that

$$(4.47) \quad \tilde{\tau}_N^2 = \operatorname{Var} \tilde{T}_N^+ = \tau_N^2 + 2 \operatorname{Cov}(T_N^+, \tilde{T}_N^+ - T_N^+) + \sigma^2(\tilde{T}_N^+ - T_N^+) = 1 + O(N^{-1/6}),$$

and that

$$(4.48) \quad \mathbb{E}(\tilde{T}_N^* - T_N^*)^2 = \frac{\sigma^2(\tilde{T}_N^+ - T_N^+) - (\tilde{\tau}_N - \tau_N)^2}{\tilde{\tau}_N \tau_N} = O(N^{-1-\delta}).$$

Because of (4.47), arguing as in Lemma 4.3, we derive that for some positive numbers  $B$ ,  $\beta$  and  $\gamma$

$$|\tilde{\psi}_N^*(t)| \leq B N^{-\beta \log N}, \quad \log N \leq |t| \leq \gamma N^{1/2},$$

which implies

$$(4.49) \quad \int_{\log N \leq |t| \leq \gamma N^{1/2}} |t|^{-1} |\tilde{\psi}_N^*(t) - e^{-t^2/2}| dt = O(N^{-1/2}).$$



Furthermore, it follows from (4.48) that uniformly in  $t$

$$|\tilde{\psi}_N^*(t) - \psi_N^*(t)| = O(|t| E|\tilde{T}_N^* - T_N^*|) = O(|t| N^{-1/2-\delta/2})$$

which, together with (4.45), implies that

$$(4.50) \quad \int_{|t| \leq \log N} |t|^{-1} |\tilde{\psi}_N^*(t) - e^{-t^2/2}| dt = O(N^{-1/2}).$$

Finally (4.1), (4.49) and (4.50) completes the proof for exact scores.

PROOF OF THEOREM 2.2. Suppose that Assumption E is satisfied, then

$$(4.51) \quad f(x-\theta) - f(x) = -\theta f'(x) + \int_x^{x-\theta} (x-\theta-t) f''(t) dt.$$

Since  $f'(x) = -f'(-x)$ , we have uniformly in  $j=1, 2, \dots, N$

$$(4.52) \quad \int J^2(F^+(|x|))(f(x-\theta_j) - f(x)) dx = O(\theta_j^2).$$

Hence it follows from (4.14), (4.24), (4.25) and (4.47) that

$$(4.53) \quad \tau_N^2 = 1 + O(N^{-1/3-\epsilon}), \quad \epsilon = \min(1/6, \delta/2)$$

for approximate scores as well as exact scores.

We now consider the asymptotic expansion of  $ET_N^+$ . Since  $E\{|e_j - e_{jj}|^{1+\delta}\{X_j\} = O(N^{-1/2-\delta/2})$ , applying Lemma 3.1, we have

$$(4.54) \quad ET_N^+ - \sum_{j=1}^N c_j EJ(e_{jj}) \operatorname{sgn} X_j = O\left(\sum_{j=1}^N |c_j \theta_j| N^{-1/2-\delta/2}\right) = O(N^{-1/2-\delta/2}).$$

Furthermore, it follows from (3.2), (4.5) and (4.17) that

$$\sum_{j=1}^N c_j EJ(e_{jj}) \operatorname{sgn} X_j = \sum_{j=1}^N c_j EJ(F^+(|X_j|)) \operatorname{sgn} X_j + O(N^{-1/2}) = \mu_N + O(N^{-1/2}),$$

which, together with (4.54) and Lemma 3.2, implies that

$$(4.55) \quad ET_N^+ = \mu_N + O(N^{-1/2})$$

for approximate scores as well as exact scores. Hence (2.2) follows from (4.53), (4.55) and Theorem 2.1. The proof follows.

## 5. Comments

Our model, the sequence of location alternatives, may not be contiguous since we have not required finite Fisher's information.

Under the same model as ours, Puri and Wu [17] have derived a bound of order  $O(N^{-1/2+\alpha})$ ,  $\alpha > 0$  for the rate of convergence with approximate scores. But we note that our assumptions on p.d.f.  $f(x)$  is much weaker than theirs while we assume Lipschitz's condition of order  $\delta$ ,  $0 < \delta \leq 1$ , on the first derivative  $J'(t)$  of the score generating function, which is more restrictive than theirs (only first Radon—

Nikodym derivative). However, assuming only the first derivative, we can derive a bound of order  $O(N^{-1/2} \log N)$  for approximate as well as exact scores even under much weaker assumptions of this paper. This is due to the more powerful method of applying the sharp bound obtained by Seoh [19] and then invoking Esséen's smoothing lemma.

Now we introduce a more general model of near alternatives, say model  $NA$ ;  $X_{N1}, X_{N2}, \dots, X_{NN}$  are independent and  $X_{Nj}$ ,  $1 \leq j \leq N$ , has a density  $f(x, \theta_{Nj}) \in \mathcal{F}$ , where  $\mathcal{F} = \{f(x, \theta): \theta \in I\}$ ,  $f(x, \theta)$  is an absolutely continuous p.d.f. and  $I$  is an open interval containing zero. Under this model, Hušková [10], [11] has established Berry—Esséen theorems for unsigned linear rank statistics with bounded score generating functions (with unbounded second derivatives). We can, of course, generalize ours to this model  $NA$ , but we have not, not only because ours is quite satisfactory in practical purposes but also because stronger assumptions are needed and doing so does not cover the important cases of normal scores.

We should note that Müller-Funk and Witting [13] derived a bound of order  $O(N^{-1/2}(\log N)^2)$  under the assumption that underlying distributions are i.i.d. but not necessarily symmetric, and that Ralescu and Puri [18] obtained a bound of order  $O(N^{-1/2+\alpha})$ ,  $\alpha > 0$ , under the assumption that the underlying distributions are only independent (not necessarily identical or symmetric). Their results include normal scores case but their statistics are special cases of ours, i.e., the case when  $c_{N1} = c_{N2} = \dots = c_{NN} = 1$ .

With regression set up as ours, it is hoped that Berry—Esséen's bound is attainable under general alternatives of Ralescu and Puri [18]. We believe that Berry—Esséen bound for normal scores be true at least under contiguous location alternatives, but so far we have not overcome computational difficulties.

#### REFERENCES

- [1] BERGSTRÖM, H. and PURI, M. L., Convergence and remainder terms in linear rank statistics, *Ann. Statist.* **5** (1977), 671—680. *MR* **56** #9772.
- [2] BICKEL, P. J., Edgeworth expansions in nonparametric statistics, *Ann. Statist.* **2** (1974), 1—20. *MR* **50** #3444.
- [3] DOES, R. J. M. M., *Higher order asymptotics for simple linear rank statistics*, Mathematisch Centrum, Amsterdam, 1982.
- [4] DUPAČ, V. and HÁJEK, J., Asymptotic normality of simple linear rank statistics under alternatives II, *Ann. Math. Statist.* **40** (1969), 1992—2017. *MR* **40** #6701.
- [5] FELLER, W., *An introduction to probability theory and its applications*, Vol. 2, 2nd edition, Wiley, New York—London—Sydney, 1971. *MR* **42** #5292.
- [6] HÁJEK, J., Asymptotically most powerful rank order tests, *Ann. Math. Statist.* **33** (1962), 1129—1147. *MR* **26** #863.
- [7] HÁJEK, J., Asymptotic normality of simple linear rank statistics under alternatives, *Ann. Math. Statist.* **39** (1968), 325—346. *MR* **36** #6037.
- [8] HÁJEK, J. and ŠIDÁK, Z., *Theory of rank tests*, Academic Press, New York—London, 1967. *MR* **37** #4925.
- [9] HUŠKOVÁ, M., Asymptotic distribution of simple linear rank statistics for testing symmetry, *Z. Wahrscheinlichkeitstheorie und Verw. Gebiete* **14** (1969—70), 308—322. *MR* **43** #2787.
- [10] HUŠKOVÁ, M., The rate of convergence of simple linear rank statistics under hypothesis and alternatives, *Ann. Statist.* **5** (1977), 658—670. *MR* **56** #16911.
- [11] HUŠKOVÁ, M., The Berry—Esséen theorem for rank statistics, *Comment. Math. Univ. Carolin.* **20** (1979), 399—415. *MR* **80k**: 62032.
- [12] JUREČKOVÁ, J. and PURI, M. L., Order of normal approximation for rank test statistics distribution, *Ann. Probability* **3** (1975), 526—533. *MR* **51** #14350.

- [13] MÜLLER-FUNK, U. and WITTING, H., On the rate of convergence in the CLT for signed linear rank statistics, *Nonparametric statistical inference*, Vol I, (B. V. Gnedenko, M. L. Puri and I. Vincze, Eds.), North-Holland Publishing Company, 1982, 637—652.
- [14] PETROV, V. V., *Sums of independent random variables*, Springer-Verlag, Berlin, 1975. MR 52 #9335.
- [15] PURI, M. L. and SEOH, M., Berry—Esséen theorems for signed linear rank statistics with regression constants, *Limit Theorems in Probability and Statistics*, (P. Révész, Ed.), North-Holland Publishing Company, 1984, 875—905.
- [16] PURI, M. L. and SEOH, M., Edgeworth expansion for signed linear rank statistics with regression constants, *J. Statist. Planning Inference* 10 (1984), 137—149.
- [17] PURI, M. L. and WU, T. J., The order of normal approximation for signed linear rank statistics, *Teor. Veroyatnost. i Primen.* 31 (1986), 156—163.
- [18] RALESCU, S. S. and PURI, M. L., On the rate of convergence in the central limit theorem for signed rank statistics, *Advances in Appl. Math.* 6 (1985), 23—51.
- [19] SEOH, M., A bound on characteristic functions of signed linear rank statistics, *Kyungpook Math. J.* 23 (1983), 1—12.
- [20] SEOH, M., RALESCU, S. S. and PURI, M. L., Cramér type large deviations for generalized linear rank statistics under alternatives, *Ann. Probability* 13 (1985) 115—125.
- [21] VAN ZWET, W. R., On the Edgeworth expansion for the simple linear rank statistics, *Nonparametric statistical inference*, Vol. II, (B. V. Gnedenko, M. L. Puri and I. Vincze, Eds.), North-Holland Publishing Company, 1982, 889—909.

(Received May 14, 1983)

DEPARTMENT OF MATHEMATICS AND STATISTICS  
WRIGHT STATE UNIVERSITY  
DAYTON, OH 45435

DEPARTMENT OF MATHEMATICS  
UNIVERSITY OF INDIANA  
SWAIN HALL EAST  
BLOOMINGTON, IN 47405  
U.S.A.



## REDUCTION OF EQUIVARIANT BORDISM GROUPS

S. S. KHARE<sup>1</sup>

### § 1. Introduction

In [1] and [2], we have seen that the fixed data of an elementary 2-group contained in the center of a group  $G$  determines  $G$ -bordism classes. The technique used in [1] is that of Stong [3]. In [2], we have used induction technique. In this note we look at this problem from an entirely different point of view. Here our effort is to reduce the  $G$ -bordism group to  $H$ -bordism group, where  $H$  is a Sylow 2-group in  $G$ . We have used the technique of restriction and extension.

### § 2. Preliminaries

Let  $G$  be a finite group and  $H$  a subgroup of  $G$ . Let  $(X, A, \psi)$  be a topological  $G$ -triple. Let  $\mathcal{J}' \subset \mathcal{J}$  be families in  $G$ . Following the notations of Stong [4], the restriction map

$$R_H^G: \mathfrak{N}_*(G; \mathcal{J}, \mathcal{J}')(X, A, \psi) \rightarrow \mathfrak{N}_*(H; \mathcal{J}_H, \mathcal{J}'_H)(X, A, \psi)$$

is defined by

$$R_H^G([M, M_0, M_1, \theta, f]) = [M, M_0, M_1, \theta|_{H \times M}, f],$$

where  $\mathcal{J}_H$  is the family in  $H$  consisting of all subgroups of  $H$  which are in  $\mathcal{J}$ .

Next let  $(\mathcal{J}, \mathcal{J}')$  be a pair of families in  $H$ . Let  $\mathcal{J}^G$  be the family of subgroups of  $G$  which are conjugate to an element of  $\mathcal{J}$ . It is easy to see that  $(\mathcal{J}^G)_H = \mathcal{J}$ . The extension map

$$E_H^G: \mathfrak{N}_*(H; \mathcal{J}, \mathcal{J}')(X, A, \psi) \rightarrow \mathfrak{N}_*(G; \mathcal{J}^G, \mathcal{J}'^G)(X, A, \psi),$$

is given by

$$E_H^G([M, M_0, M_1, \theta, f]) = \left[ \frac{G \times M}{H}, \frac{G \times M_0}{H}, \frac{G \times M_1}{H}, \theta', f' \right],$$

where the  $H$ -action on  $G \times M$  is given as  $h(g, x) = (gh^{-1}, hx)$ , the action  $\theta'$  of  $G$

on  $\frac{G \times M}{H}$  is given as  $\theta'(g', \overline{(g, x)}) = \overline{(g' \cdot g, x)}$  and  $f': \frac{G \times M}{H} \rightarrow X$  is given by

$f'(\overline{(g, x)}) = \overline{gf(x), (g, x)}$  being the orbit of  $(g, x)$  in  $\frac{G \times M}{H}$ .

<sup>1</sup> The author was partially supported by D.A.E. grant.

1980 *Mathematics Subject Classification*. Primary 57R05.

*Key words and phrases*. Equivariant bordism, equivariant trivial normal bundle,  $G$ -representation, computable equivariant bordism group.

### § 3. Reduction of groups

Let  $H$  be a normal subgroup of  $G$  and  $\{g_1, \dots, g_m\}$  be a set of representatives of the cosets of  $H$  in  $G$ . Let  $M_i$  be  $M$  with the action  $\theta_i$  of  $H$  on  $M$  given by  $\theta_i(h, x) = \theta(g_i^{-1} h g_i, x)$  and  $f_i: M_i \rightarrow X$  be given by  $f_i(x) = g_i f(x)$ .

THEOREM 3.1. Let  $[M, M_0, M_1, \theta, f] \in \mathfrak{N}_*(H; \mathcal{J}, \mathcal{J}')(X, A, \psi)$ . Then

$$R_H^G E_G^H([M, M_0, M_1, \theta, f]) = \sum_{i=1}^m [M_i, (M_0)_i, (M_1)_i, \theta_i, f_i].$$

PROOF. Consider the manifold  $\frac{G}{H} \times M$  with the action  $\alpha$  of  $H$  given as

$$\alpha(h, (g_i H, x)) = (g_i H, (g_i^{-1} h g_i) x).$$

Consider the map  $\beta: \frac{G \times M}{H} \rightarrow \frac{G}{H} \times M$  defined as  $\beta((\overline{g}, x)) = (g_i H, h x)$ , where  $g = g_i h$  for some unique coset representative  $g_i$  and  $h \in H$ . It is easy to verify that  $\beta$  is an equivariant diffeomorphism. Let  $f'': \frac{G}{H} \times M \rightarrow X$  be the map defined by  $f''(g_i H, x) = g_i f(x)$ . We have the following commutative diagram.

$$\begin{array}{ccc} \frac{G \times M}{H} & \xrightarrow{\beta} & \frac{G}{H} \times M \\ & \searrow f' & \swarrow f'' \\ & X & \end{array}$$

Therefore we have

$$\begin{aligned} R_H^G E_G^H([M, M_0, M_1, \theta, f]) &= \left[ \frac{G \times M}{H}, \frac{G \times M_0}{H}, \frac{G \times M_1}{H}, \theta' | H \times \left( \frac{G \times M}{H} \right), f' \right] \\ &= \left[ \frac{G}{H} \times M, \frac{G}{H} \times M_0, \frac{G}{H} \times M_1, \alpha, f'' \right] = \\ &= \sum_{i=1}^m [M_i, (M_0)_i, (M_1)_i, \theta_i, f_i]. \end{aligned}$$

COROLLARY 3.2. If  $H$  is contained in the center, then

$$R_H^G E_G^H([M, M_0, M_1, \theta, f]) = \sum_{i=1}^m [M, M_0, M_1, \theta, f_i].$$

COROLLARY 3.3. Let  $H$  be contained in the center of  $G$ . Let there be coset representatives  $g_1, \dots, g_m$  of  $G/H$  which act trivially on  $X$  up to homotopy. Then  $R_H^G E_G^H$  is the identity map, if  $G/H$  is of odd order, and is the zero homomorphism, if  $G/H$  is of even order.



PROOF. Since  $g_1, \dots, g_m$  act trivially on  $X$  up to homotopy,  $[M, M_0, M_1, \theta, f_i] = [M, M_0, M_1, \theta, f], \forall i$ . This proves the Corollary.

THEOREM 3.4. Let  $G$  be a finite group with nontrivial Sylow 2-group  $H$  contained in the center of  $G$ . Let there be coset representatives  $g_1, \dots, g_m$  of  $G/H$  which act on  $X$  trivially up to homotopy. Suppose  $(\mathcal{J}, \mathcal{J}')$  is a pair of families of 2-groups in  $G$ . Then

$$\mathfrak{N}_*(G; \mathcal{J}, \mathcal{J}')(X, A, \psi) \overset{R_H^G}{\approx} \mathfrak{N}_*(H; \mathcal{J}_H, \mathcal{J}'_H)(X, A, \psi).$$

PROOF. From the Proposition 7.1 of [4], one gets that  $E_G^H R_H^G$  is an isomorphism. Since  $G/H$  is of odd order, this together with the Corollary 3.3 gives the result.

REMARK 3.5. Theorem 3.4 reduces the problem of singular  $G$ -bordism theory completely to the problem of  $H$ -bordism theory, where  $H$  is Sylow 2-group of  $G$  contained in the center of  $G$ .

#### § 4. Applications

We first apply our results to the computability problem of equivariant bordism groups. The equivariant bordism groups  $\mathfrak{N}_*(G; \mathcal{J}, \mathcal{J}')(X, A, \psi)$  are said to be *computable*, if they are naturally isomorphic to a direct sum of ordinary unoriented bordism groups with possibly dimension shifts.

PROPOSITION 4.1. Let  $G$  be a finite group with nontrivial Sylow 2-group  $H$  in the center of  $G$ . Let there be representatives  $g_i$  of the cosets of  $H$  in  $G$  which act on  $X$  trivially up to homotopy. Suppose  $\mathfrak{A}_2$  is the family of all 2-groups of  $G$ . Then  $\mathfrak{N}_*(G; \mathfrak{A}_2)(X, A, \psi)$  is computable.

PROOF. By the Theorem 3.4, one gets

$$\mathfrak{N}_*(G, \mathfrak{A}_2)(X, A, \psi) \overset{R_H^G}{\approx} \mathfrak{N}_*(H; \mathfrak{A})(X, A, \psi),$$

where  $\mathfrak{A}$  is the family of all subgroups of  $H$ . By the Proposition 9.2 of [4],  $\mathfrak{N}_*(H; \mathfrak{A})(X, A, \psi)$  is computable.

Next we apply our results to show that the triviality of the equivariant normal bundle of the fixed point set in the manifold under the elementary 2-group contained in the center of the group  $G$  determines the  $G$ -bordism classes.

Let  $[M^n, \theta] \in \mathfrak{N}_*(G; \mathfrak{A}_2)$ . Let  $\mathcal{P}_2$  be the family of all proper subgroups of Sylow 2-groups of the group. Let  $\mathbf{Z}_2^\otimes$  be the nontrivial subgroup of all elements of order 2 contained in the center together with the identity element. Let  $F$  be the fixed point set of  $\mathbf{Z}_2^\otimes$  in  $M^n$ . Consider  $F = \bigcup_{l=0}^n F^l$ , where  $F^l$  is the  $l$ -dimensional component of  $F$ . Let  $v_l$  be the normal bundle of  $F^l$  in  $M^n$  and  $D(v_l)$  be the disc bundle of  $v_l$ . Let  $\theta_l$  be the action on  $D(v_l)$  induced from the action  $\theta$  on  $M$ . Let  $H$  be the Sylow 2-group contained in the center of  $G$ .

DEFINITION 4.2.  $F$  is said to have equivariant trivial normal bundle in  $M^n$  in restricted sense, if  $H/\mathbf{Z}_2^\otimes$  acts trivially on  $F$  and  $\exists$  positive dimensional  $G$ -repre-

sentations  $(W_l, \varphi_l)$ ,  $1 \leq l \leq n$ , such that in  $\mathfrak{N}_*(G; \mathfrak{U}_2, \mathcal{P}_2)$

$$[D(v_l), \theta_l] = [F^l][D(W_l, \varphi_l)],$$

where  $D(W_l)$  is the unit disc of  $W_l$ .

**THEOREM 4.3.** *Let  $[M^n, \theta] \in \mathfrak{N}_n(G; \mathfrak{U}_2)$ . If  $F$  has equivariant trivial normal bundle in  $M^n$  in restricted sense, then  $[M^n, \theta]$  is zero in  $\mathfrak{N}_n(G; \mathfrak{U}_2)$ .*

**PROOF.** From the Theorem 3.4, one has

$$\mathfrak{N}_n(G; \mathfrak{U}_2, \mathcal{P}_2) \overset{R_H^G}{\approx} \mathfrak{N}_n(H; \mathfrak{U}, \mathcal{P}),$$

where  $\mathcal{P}$  is the family of all subgroups of  $H$  other than  $H$  itself. Therefore  $F$  will have equivariant trivial normal bundle in  $M^n$  in the sense of Definition 4.2. Hence by the Theorem 4.3 of [1], one gets that  $[M^n, \theta]$  is zero in  $\mathfrak{N}_n(H; \mathfrak{U})$ . Further one also has

$$\mathfrak{N}_n(G; \mathfrak{U}_2) \overset{R_H^G}{\approx} \mathfrak{N}_n(H; \mathfrak{U}).$$

Therefore  $[M^n, \theta]$  is zero in  $\mathfrak{N}_n(G; \mathfrak{U}_2)$ . This completes the proof of the Theorem.

#### REFERENCES

- [1] KHARE, S. S.,  $(\mathcal{J}, \mathcal{J}')$ -free bordism, characteristic numbers and stationary point set, *Acta Math. Hungar.* 45 (1985), 45—52.
- [2] KHARE, S. S., Compact Lie group action and equivariant bordism, *Proc. Amer. Math. Soc.* 92 (1984), 297—300.
- [3] STONG, R. E., Equivariant bordism and  $(\mathbb{Z}_2)^k$ -action, *Duke Math. J.* 37 (1970), 779—785. *MR* 42#6847.
- [4] STONG, R. E., *Unoriented bordism and actions of finite groups*, A.M.S. Memoir, No. 103, 1970. *MR* 42#8522.

(Received May 23, 1983)

DEPARTMENT OF MATHEMATICS  
NORTH-EASTERN HILL UNIVERSITY  
BIJNI COMPLEX  
SHILLONG 793 003  
MEGHALAYA  
INDIA

# PACKING OF HOMOTHETIC DISCS OF $n$ DIFFERENT SIZES

L. FEJES TÓTH

We shall prove the following

**THEOREM.** *We consider a packing of the plane with homothetic copies of a convex disc  $c$  of at most  $n$  different sizes. If  $d_n(c)$  is the supremum of the upper densities of all such packings then*

$$(1) \quad d_n(c) \cong \frac{3^n - 1}{3^n}$$

*and equality holds iff  $c$  is a triangle.*

The main interest of this theorem lies in the various unsolved problems connected with it which we shall discuss later.

The inequality (1) is a simple consequence of a theorem of Fáry [1]: If  $d(c)$  is the density of the densest lattice-packing of translates of  $c$  then  $d(c) \cong 2/3$ , and equality holds iff  $c$  is a triangle.

Let  $P_\varepsilon$  be a densest lattice-packing of translates of  $\varepsilon c$ . We define the packing  $P^n$  inductively as follows. Let be  $P^0 = P_1$ . Let  $P^n$  be the packing consisting of the discs of  $P^{n-1}$ , and those discs of  $P_\varepsilon^n$  which are disjoint of the discs of  $P^{n-1}$ . Let  $p_n$  be the part of the plane not covered by  $P^n$ . Since the density of  $P^0$  is equal to  $d = d(c)$ , the density of  $p_0$  is equal to  $1 - d$ . We now suppose that  $\varepsilon \ll 1$ . Then, roughly speaking, in  $P^1$   $p_0$  is filled by translates of  $\varepsilon c$  with density  $d$ , so that the density of  $P^1$  is equal to  $(1 - d)^2$ , and so on. Thus, for sufficiently small  $\varepsilon$ , the density of  $P^n$  will be arbitrarily close to  $1 - (1 - d)^n$ , showing that

$$d_n(c) \cong 1 - (1 - d)^n.$$

By Fáry's theorem this implies (1) with strict inequality for any disc other than a triangle.

We still have to show that for a triangle  $\Delta$  in (1) equality holds. We shall show this by proving that

$$(2) \quad d_n(\Delta) \cong \frac{3^n - 1}{3^n}.$$

Let  $\Pi_n = \{\Delta_i\} = \{A_i B_i C_i\}$  be a packing of homothetic images of  $\Delta = ABC$  of  $n$  different sizes. Let  $\Delta$  be in a position such that a half-line issuing from  $C$  and

1980 *Mathematics Subject Classification*. Primary 52A45.

*Key words and phrases*. Packing, covering, density.

intersecting the side  $AB$  points vertically upwards. We associate with  $\Delta_i$  a region  $r_i$  consisting of the points of  $\Delta_i$  and those points above  $\Delta_i$  which do not lie in or above another triangle of  $\Pi_n$ . We denote a domain and its area with the same symbol and claim that if  $\Delta_i$  is a smallest triangle in  $\Pi_n$  then

$$(3) \quad \Delta_i \leq \frac{2}{3} r_i.$$

To see this reflect  $\Delta_i$  in the midpoint  $M$  of the side  $A_i B_i$  obtaining the triangle  $\bar{\Delta}$ . If no triangle  $\Delta_j$  overlaps  $\bar{\Delta}$  then  $r_i \supset \Delta_i \cup \bar{\Delta}$ , and consequently  $r_i \geq 2\Delta_i$ . We now suppose that  $\Delta_j \cap \bar{\Delta} > 0$ . Since no translate of  $\lambda\Delta_i$  with  $\lambda \geq 1$  can overlap the region  $w = \bar{\Delta} \setminus \Delta_j$  without overlapping  $\Delta_i$  or  $\Delta_j$ , we have  $r_i \supset \Delta_i \cup w$ , and consequently  $r_i \geq \Delta_i + w$ . It is easy to see that, for various admissible positions of  $\Delta_j$ ,  $w$  attains its minimum if  $C_j$  coincides with  $M$ , so that  $w$  consists of two congruent triangles. Now  $w = \bar{\Delta}/2 = \Delta_i/2$ , so that, in accordance with our statement,  $r_i \geq 3\Delta_i/2$ .

Let  $\pi_n$  be the part of the plane not covered by  $\Pi_n$ , and let  $\delta_n$  be the density of  $\pi_n$ . Since the regions  $r_i$  do not overlap, the inequalities (3) imply that the density of  $\Pi_1$  is at most  $2/3$ , so that  $\delta_1 \geq 1/3$ . Now we make the inductive assumption that  $\delta_{n-1} \geq 1/3^{n-1}$ . Since the regions  $r_i$  belonging to the smallest triangles in  $\Pi_n$  are all contained in  $\pi_{n-1}$  but do not fill it completely, we have, because of (3), for  $n > 1$ ,  $\delta_n > 1/3$   $\delta_{n-1} \geq 1/3^n$ . Thus the density of  $\Pi_n$  is  $\leq 1 - 3^{-n}$ , and equality is claimed only for  $n=1$ . This completes, along with the proof of (2), the proof of the theorem.

The inequality  $d_1(\Delta) \leq 2/3$  is a special case of a general theorem [2] according to which the density of a packing of translates of a convex disc never exceeds the density of the densest lattice-packing of these discs. The proof of (2), due to A. Bezdek, is based upon a new proof [3] of this theorem.

**CONJECTURE 1.** *If  $d=d(c)$  is the density of the densest lattice-packing of translates of a convex disc  $c$  then*

$$(4) \quad d_n(c) = 1 - (1-d)^n.$$

Besides triangles (4) holds also for circles [4], and consequently for ellipses. By a slight modification of the proof for triangles (cf. [3]) the equality (4) can be proved for a more general class of (not necessarily convex) discs defined as follows. Divide a parallelogram by a simple Jordan-arc connecting opposite vertices in two parts, and consider any of these parts. A simple example is given by a circular sector of angle  $\leq 90^\circ$ .

It would be nice to complete the above theorem by a dual counterpart which we phrase as

**CONJECTURE 2.** *We consider a covering of the plane with homothetic copies of a convex disc  $c$  of at most  $n$  different sizes. If  $D_n(c)$  is the infimum of the lower densities of all such coverings then*

$$(5) \quad D_n(c) \leq \frac{3^n}{3^n - 1}$$

*and equality holds iff  $c$  is a triangle.*

Fáry's above-mentioned theorem has a dual counterpart [1]: If  $D(c)$  is the density of the thinnest lattice-covering by translates of  $c$  then  $D(c) \leq 3/2$ , and equality holds iff  $c$  is a triangle. Thus

$$(6) \quad D_1(c) \leq 3/2$$

and equality cannot hold but for triangles. It is known [5] that the density of a covering with translates of a centro-symmetric convex disc is never less than the density of the thinnest lattice-covering. Unfortunately we do not know whether this holds without the condition of central symmetry. Therefore the question whether for a triangle we have actually equality in (6) remains outstanding.

We claim that if  $\Delta$  is a triangle then

$$(7) \quad D_n = D_n(\Delta) \leq \frac{3^n}{3^n - 1}.$$

The proof of (7) was proposed in 1982 as a problem at the Miklós Schweitzer Mathematical Competition of the János Bolyai Mathematical Society. For the sake of completeness we present the proof.

We start with a regular hexagonal tiling. Circumscribing about the hexagons similarly situated regular triangles we obtain a lattice-covering by triangles of density  $3/2$ . Let the inradius of the triangles be 1. Translate the sides of each triangle inwards through a distance  $\varrho \leq 1/3$  obtaining a lattice  $L_\varrho$  of triangles of inradius  $1 - \varrho$  with oppositely situated triangular gaps of inradius  $\varrho$ . The density of the triangles is  $\frac{3}{2}(1 - \varrho)^2$ . Since there are twice as many gaps as triangles, the density of the gaps is equal to  $2 \cdot \frac{3}{2} \varrho^2 = 3\varrho^2$ .

Now we cover each gap with some kinds of discs of total area  $t$  times the area of a gap. We consider the joint density

$$D = \frac{3}{2}(1 - \varrho)^2 + 3t\varrho^2$$

of the triangles and the discs, and observe that  $D$  has a minimum,  $3t/(2t+1)$ , at  $\varrho = 1/(2t+1)$ .

We construct the lattice  $L_\varrho$  with  $\varrho = 1/\left(2 \cdot \frac{3}{2} + 1\right) = 1/4$ , and cover the gaps with "very small" congruent triangles homothetic to those of  $L_\varrho$  with density close to  $3/2$ . In this way we get a covering with two kinds of homothetic triangles with a density arbitrarily close to  $3 \cdot \frac{3}{2} / \left(2 \cdot \frac{3}{2} + 1\right) = 9/8$ , showing that  $D_2 \leq 9/8$ . Starting with  $D_1 \leq 3/2$  rather than with  $3/2$ , we see that  $D_2 \leq 3D_1/(2D_1+1)$ . In the next step we start with  $\varrho = 1/\left(2 \cdot \frac{9}{8} + 1\right)$  and cover the gaps in  $L_\varrho$  with two kinds of homothetic triangles with density close to  $9/8$ , showing that  $D_3 \leq 3D_2/(2D_2+1)$ , and so on, or explicitly  $D_n \leq 3^n/(3^n - 1)$ , as stated.

In an arrangement of discs we call the density of the part of the plane covered by the discs *covering measure* of the arrangement. Let  $x = \frac{3}{2}(1-\varrho)^2$  be the density, and  $m = 1 - 3\varrho^2$  the covering measure of  $L_\varrho$ . Eliminating  $\varrho$  we obtain  $m$  in terms of  $x$ :

$$m(x) = 2(\sqrt{6x} - x - 1), \quad 2/3 \leq x \leq 3/2.$$

We extend the definition of the function  $m(x)$  to any  $x \geq 0$  by  $m(x) = x$  for  $0 \leq x < 2/3$  and  $m(x) = 1$  for  $x > 3/2$ . The graph of  $m(x)$  is exhibited in Fig. 1.

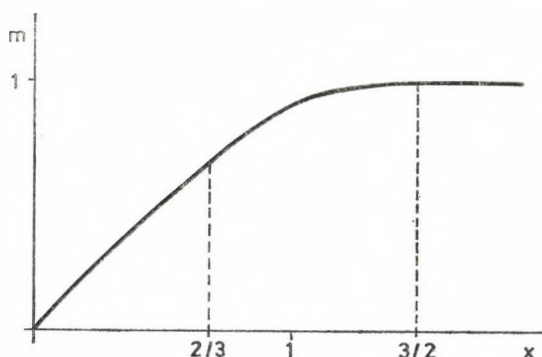


Fig. 1

**CONJECTURE 3.** *We consider a lattice of translates of a convex disc  $c$  of prescribed density  $x$ . Let  $M(x; c)$  be the maximum of the covering measure of all such lattices. Then  $M(x; c) \leq m(x)$  and for  $2/3 \leq x \leq 3/2$  equality holds iff  $c$  is a triangle.*

By the construction used in the proof of (7) this conjecture would imply (5) with strict inequality for any disc other than a triangle. Moreover, Conjecture 3 offers a possibility of uniting Fáry's packing and covering theorem in one general theorem.

An analogue of (4) which would yield  $D_n(c)$  in terms of  $D(c)$  is not to be expected because  $D_n(c)$  depends on many other parameters. Let  $\bar{M}(x; c)$  be the supremum of the covering measure of translates of  $c$  scattered in the plane with density  $x$ . It is very likely that an algorithm similar to the proof of (5) based on  $\bar{M}(x; c)$  would yield the value of  $D_n(c)$ . G. Fejes Tóth [6] constructed a convex disc  $c$  such that for some  $x$  we have  $\bar{M}(x; c) > M(x; c)$ . He showed [7] that if  $c$  is centro-symmetric then  $\bar{M}(x; c)$  is the least concave function not less than  $M(x; c)$ . We do not know whether this is true for any convex disc. By rather complicated arguments G. Fejes Tóth succeeded in determining  $D_2(c)$  for a circle.

## REFERENCES

- [1] FÁRY, I., Sur la densité des réseaux de domaines convexes, *Bull. Soc. Math. France* **78** (1950), 152—161. *MR* **12**—526.
- [2] ROGERS, C. A., The closest packing of convex two-dimensional domains, *Acta Math.* **86** (1951) 309—321. *MR* **13**—768.
- [3] FEJES TÓTH, L., On the densest packing of convex discs, *Mathematika* **30** (1983), 1—3.
- [4] FEJES TÓTH, L., Ausfüllung eines konvexen Bereiches durch Kreise, *Publ. Math. Debrecen* **1** (1949), 92—94. *MR* **12**—123.
- [5] BAMBAH, R. P. and ROGERS, C. A., Covering the plane with convex sets, *J. London Math. Soc.* **27** (1952), 304—314. *MR* **13**—971.
- [6] FEJES TÓTH, G., On the intersection of a convex disc and a polygon, *Acta Math. Acad. Sci. Hungar.* **29**(1977), 149—153. *MR* **57**#1273.
- [7] FEJES TÓTH, G., Covering the plane by convex discs, *Acta Math. Acad. Sci. Hungar.* **23** (1972), 263—274. *MR* **47**#9431.

(Received June 7, 1983)

MTA MATEMATIKAI KUTATÓ INTÉZETE  
P.O. BOX 127  
H—1364 BUDAPEST  
HUNGARY





# A REMARK ON MATHUR'S PAPER

## SIMPLE PROOFS OF TELYAKOVSKII—GOPENGAUZ'S THEOREM

K. B. SRIVASTAVA

### 1. Preliminaries and results

In his paper K. K. Mathur [3] established Jackson's theorem for functions  $f \in C[-1, 1]$  by employing Freud's [1] interpolation polynomials  $S_n(f, x)$  of degree  $4n$  at most constructed on the nodes

$$(1.1) \quad x_{kn} = \cos \frac{2k-1}{2n+1} \pi \quad k = 1, \dots, n; \quad n = 1, 2, \dots \quad (1)$$

defined by

$$(1.2) \quad S_n(f, x) = -xf(-1) + \sum_{k=1}^n [f(x_k) + xf(-1)] \mu_k(x)$$

where

$$(1.3) \quad \mu_k(x) = \left( \frac{1+x}{1+x_k} \right)^2 [v_k(x) \varphi_k^4(x) + 2(1+x_k)(x-x_k) \varphi_k^3(x) X_{n-1}(x_k, x)]$$

$$(1.4) \quad v_k(x) = 1 + \frac{1-2x_k}{1-x_k^2} (x-x_k)$$

$$(1.5) \quad x_{n-1}(x_k, x) = \frac{2}{n+\frac{1}{2}} \left[ \sum_{r=1}^{n-1} P_r^{((-1/2), (1/2))}(x_k) P_r^{((-1/2), (1/2))}(x) \right] \quad (2)$$

and

$$(1.6) \quad \begin{aligned} \varphi_k(x) &= \frac{P_n(x)}{P'_n(x_k)(x-x_k)} = \\ &= \frac{(-1)^{k-1} \sin \theta_k \cos \theta_{k/2} \cos \left( n + \frac{1}{2} \right) \theta}{\left( n + \frac{1}{2} \right) \cos \theta/2 (\cos \theta - \cos \theta_k)}, \quad x = \cos \theta \end{aligned}$$

being the fundamental Lagrange interpolation polynomial built on the abscissas (1.1) that are roots of

$$(1.7) \quad P_n(x) = \frac{\cos(n+1/2)\theta}{\cos \theta/2}, \quad x = \cos \theta.$$

<sup>1</sup> We shall henceforth be writing  $k$  only instead of  $kn$ .

<sup>2</sup> Throughout this note we shall denote  $P_n(x)$  to mean  $P_n^{(-\frac{1}{2}, \frac{1}{2})}(x)$  unless otherwise stated 1980 *Mathematics Subject Classification*. Primary 41A15; Secondary 41A18.

*Key words and phrases*. Interpolation polynomial, Telyakovskii-Gopengauz's theorem.

His method of proof involves an erroneous use of the properties of modulus of continuity of functions which can be made possible only when  $f(1) = -f(-1)$ . But this is not true of all continuous functions on  $[-1, 1]$ . Thus he proves his theorem for another class of functions contained in  $C[-1, 1]$  without any prior mention of the class itself.

It is worthwhile to indicate that similar methods have extensively been employed by many authors, e.g. Kis—Vértesi [11] Freud and Vértesi [2], Saxena [5, 6, 7], Sallay [8], Srivastava [9] including the original result by Freud [1] himself to reproduce Jackson's, Timan's and Telyakovskii's [10] theorems for continuous functions.

One of the main aims in this note is to correct the form of the polynomial  $S_n(f, x)$  in (1.2) by defining another interpolation process  $R_n(f; x)$  built on the abscissas (1.1) given by

$$(1.8) \quad R_n(f, x) = t(x) + \sum_{k=1}^n [f(x_k) - t(x)] \mu_k(x)$$

where  $t(x)$  is a linear function defined by

$$(1.9) \quad t(x) = \frac{1+x}{2} f(1) + \frac{1-x}{2} f(-1).$$

We now do not require to restrict the class of functions. With the help of  $R_n(f, x)$  we shall establish Timan's theorem. Namely, we prove the following

**THEOREM 1.** *Let  $f \in C[-1, 1]$ , then, for the sequence of interpolation polynomial  $R_n(f, x)$  given by (1.8), we have, uniformly in  $x \in [-1, 1]$*

$$(1.10) \quad |f(x) - R_n(f, x)| \leq C_1 \omega_f \left( \frac{\sqrt{1+x}}{n} + \frac{1}{n^2} \right)$$

where  $\omega_f(\cdot)$  is the usual modulus of continuity of  $f$ .

Using different techniques we can state

**THEOREM 2.** *For  $f \in C[-1, 1]$ , we have the inequality*

$$(1.11) \quad |f(x) - R_n(f, x)| \leq C_2 \omega_f \left( \frac{\sqrt{1+x}}{n} \right).$$

In order to obtain two-point Telyakovskii—Gopengauz's theorem, we modify our operators in the following manner:

$$(1.12) \quad R_n^*(f, x) = R_n(f, x) + \frac{1+x}{2} [f(1) - R_n(f, 1)].$$

It is readily verified that

$$R_n^*(f, \pm 1) = f(\pm 1).$$

With the help of  $R_n^*(f, x)$  we prove

**THEOREM 3.** *Let  $f \in C[-1, 1]$ , then, for every  $x \in [-1, 1]$ , we have uniformly*

$$(1.13) \quad |R_n^*(f, x) - f(x)| \leq C_3 \omega_f \left( \frac{\sqrt{1-x^2}}{n} \right).$$

The process  $Q_n(f, x)$  similar to (1.8) built on the abscissas

$$(1.14) \quad x_k^* = \cos \frac{2k}{2n+1} \pi, \quad k = 1, \dots, n$$

can be shown to satisfy Theorems 1 and 2 where  $x$  is replaced by  $-x$ . Then the natural modification  $Q_n^*(f, x)$  of  $Q_n(f, x)$  in the manner that follows

$$(1.15) \quad Q_n^*(f, x) = Q_n(f, x) + \frac{1-x}{2} [f(-1) - Q_n(f, -1)]$$

will again give rise to Telyakovskii—Gopengauz theorem. Still another way to get it is to define the combination of  $R_n(f, x)$  and  $Q_n(f, x)$  in the following fashion

$$(1.16) \quad W_n(f, x) = \frac{1-x}{2} R_n(f, x) + \frac{1+x}{2} Q_n(f, x).$$

We shall see how easily the Theorem 3 can be shown to hold for  $W_n(f, x)$  also.

We restrict ourselves to proving only Theorems 1, 2 and 3 and mention that exactly the same methods can be applied to obtain the corresponding theorems for  $Q_n(f, x)$  also. In order to establish our assertions we require a number of auxiliary results which we state in the following in the form of lemmas.

## 2. Auxiliary lemmas

LEMMA 1. *By setting*

$$(2.1) \quad Y_{n-1}(u, v) = \frac{2}{\left(n + \frac{1}{2}\right)} \left[1 + \sum_{r=1}^{n-1} P_r(u) P_r(v)\right]$$

we have,

$$(2.2) \quad \sum_{k=1}^n \mu_k(x) = [(1+x)Y_{n-1}(x, x)]^2.$$

PROOF. To prove (2.2), we make use of the well-known Christoffel—Darboux formula for  $P_n(x)$  and get

$$(2.3) \quad \varphi_k(x) = \frac{2(1+x_k)}{\left(n + \frac{1}{2}\right)} \left[1 + \sum_{r=1}^{n-1} P_r(x_k) P_r(x)\right] = (1+x_k)Y_{n-1}(x_k, x).$$

Now, let

$$\psi_{2n}(x) = [(1+x)Y_{n-1}(x, \xi)]^2$$

be a polynomial of degree  $\leq 2n$  satisfying the properties

$$\psi_{2n}(x_k) = \varphi_k^2(\xi), \quad \psi_{2n}(-1) = 0$$

and

$$\psi'_{2n}(x_k) = \frac{2\varphi_k^2(\xi)}{1+x_k} + 2(1+x_k)\varphi_k^2(\xi)X_{n-1}(x_k, \xi)$$

$$\psi'_{2n}(-1) = 0.$$

Then, constructing the well-known Hermite polynomial of degree  $\leq 2n+1$  on the nodes (1.1) for  $\psi_{2n}(x)$ , we readily obtain (2.2) by putting  $x=\xi$  (for details see [3]).

LEMMA 2. For  $|x| \leq 1$ , we have

$$\left|1 - \sum_{k=1}^n \mu_k(x)\right| < 3$$

and

$$\sqrt{(1-x^2)} \left|1 - \sum_{k=1}^n \mu_k(x)\right| \leq \frac{3}{n+\frac{1}{2}}.$$

PROOF. From (2.1), we have by making use of summation formulae for trigonometric polynomials

$$(1+x)Y_{n-1}(x, x) = 1 + \frac{\cos\left(n+\frac{1}{2}\right)\theta \sin\left(n-\frac{1}{2}\right)\theta}{\left(n+\frac{1}{2}\right)\sin\theta}, \quad x = \cos\theta$$

which, in turn, gives

$$(2.4) \quad \begin{aligned} 1 - [(1+x)Y_{n-1}(x, x)]^2 &= -\frac{\cos\left(n+\frac{1}{2}\right)\theta \sin\left(n-\frac{1}{2}\right)\theta}{\left(n+\frac{1}{2}\right)\sin\theta} \times \\ &\times \left[2 + \frac{\cos\left(n+\frac{1}{2}\right)\theta \sin\left(n-\frac{1}{2}\right)\theta}{\left(n+\frac{1}{2}\right)\sin\theta}\right]. \end{aligned}$$

Making use of  $|\sin n\theta| \leq n|\sin\theta|$  and  $|\cos n\theta| \leq 1$ , (2.4) gives

$$\left|1 - \sum_{k=1}^n \mu_k(x)\right| \leq \frac{n-\frac{1}{2}}{n+\frac{1}{2}} \left[ \frac{\left(n-\frac{1}{2}\right)}{\left(n+\frac{1}{2}\right)} \right] \leq 3$$

and

$$\sqrt{1-x^2} \left| 1 - \sum_{k=1}^n \mu_k(x) \right| \leq \frac{1}{n + \frac{1}{2}} \left[ 2 + \frac{n - \frac{1}{2}}{n + \frac{1}{2}} \right] < \frac{3}{n + \frac{1}{2}}$$

which proves the lemma.

For our purposes we need a more precise estimate for  $\mu_k(x)$  than in [3]. To this end, we prove the following equality.

LEMMA 3.

$$(2.5) \quad \sqrt{\frac{1+x}{2}} X_{n-1}(x_k, x) = (-1)^{k-1} \frac{\cos\left(n + \frac{1}{2}\right) \theta \sin \theta_k}{2(2n+1) \cos^3 \theta_{k/2} (\cos \theta - \cos \theta_k)} +$$

$$+ \frac{(-1)^{k-1} \sin\left(n + \frac{1}{2}\right) \theta \sin \theta}{\sin \theta_k \cos \theta_{k/2} (\cos \theta - \cos \theta_k)} + \frac{(-1)^k \cos\left(n + \frac{1}{2}\right) \theta}{\sin \theta_k \cos \theta_{k/2}} +$$

$$+ \frac{(-1)^{k-1} \cos\left(n + \frac{1}{2}\right) \theta}{2(2n+1) \sin \theta_k \cos \theta_{k/2}} \left[ \operatorname{cosec}^2 \frac{\theta + \theta_k}{2} + \operatorname{cosec}^2 \frac{\theta - \theta_k}{2} \right].$$

PROOF. We have from the definition,

$$\sqrt{\frac{1+x}{2}} X_{n-1}(x_k, x) = \frac{2}{n + \frac{1}{2}} \sum_{r=1}^{n-1} \frac{\cos\left(r + \frac{1}{2}\right) \theta}{(-\sin \theta_k)} \left[ \frac{-\left(r + \frac{1}{2}\right) \sin\left(r + \frac{1}{2}\right) \theta_k}{\cos \theta_{k/2}} + \right.$$

$$\left. + \frac{\cos\left(r + \frac{1}{2}\right) \theta_k \sin \theta_{k/2}}{2 \cos^2 \theta_{k/2}} \right] =$$

$$= \frac{2}{(2n+1) \sin \theta_k \cos \theta_{k/2}} \sum_{r=1}^{n-1} \left(r - \frac{1}{2}\right) \left\{ \sin\left(r + \frac{1}{2}\right) (\theta + \theta_k) - \sin\left(r + \frac{1}{2}\right) (\theta - \theta_k) - \right.$$

$$\left. - \frac{\sin \theta_{k/2}}{(2n+1) \sin \theta_k (\cos \theta_{k/2})} \sum_{r=1}^{n-1} \left\{ \cos\left(r + \frac{1}{2}\right) (\theta + \theta_k) + \cos\left(r + \frac{1}{2}\right) (\theta - \theta_k) \right\} = \right.$$

$$= -\frac{1}{2(2n+1) \cos^3 \theta_{k/2}} \left[ \frac{\sin n(\theta + \theta_k)}{2 \sin\left(\frac{\theta + \theta_k}{2}\right)} + \frac{\sin n(\theta - \theta_k)}{2 \sin\left(\frac{\theta - \theta_k}{2}\right)} - \right.$$

$$\left. - \cos \frac{1}{2} (\theta + \theta_k) - \cos \frac{1}{2} (\theta - \theta_k) \right] +$$

$$\begin{aligned}
& + \frac{2}{(2n+1) \sin \theta_k \cos \theta_{k/2}} \left[ \frac{n \cos n(\theta - \theta_k)}{2 \sin \frac{\theta - \theta_k}{2}} - \frac{n \cos n(\theta + \theta_k)}{2 \sin \frac{\theta + \theta_k}{2}} \right. \\
& + \frac{1}{2} \left\{ \sin \frac{\theta - \theta_k}{2} - \sin \frac{\theta + \theta_k}{2} \right\} + \frac{\sin n(\theta + \theta_k) \cos \frac{1}{2}(\theta + \theta_k)}{4 \sin^2 \frac{\theta + \theta_k}{2}} - \\
& \left. - \frac{\sin n(\theta - \theta_k) \cos \frac{1}{2}(\theta - \theta_k)}{4 \sin^2 \frac{\theta - \theta_k}{2}} \right].
\end{aligned}$$

It is easy to see that

$$\begin{aligned}
\sin n(\theta \pm \theta_k) &= \pm (-1)^{k-1} \left[ \cos \left( n + \frac{1}{2} \right) \theta \cos \frac{1}{2}(\theta \pm \theta_k) + \right. \\
&\quad \left. + \sin \left( n + \frac{1}{2} \right) \theta \sin \frac{1}{2}(\theta \pm \theta_k) \right]
\end{aligned}$$

and

$$\begin{aligned}
\cos n(\theta \pm \theta_k) &= \pm (-1)^{k-1} \left[ \sin \left( n + \frac{1}{2} \right) \theta \cos \frac{1}{2}(\theta \pm \theta_k) - \right. \\
&\quad \left. - \cos \left( n + \frac{1}{2} \right) \theta \sin \frac{1}{2}(\theta \pm \theta_k) \right].
\end{aligned}$$

Keeping these identities in mind, we have,

$$\begin{aligned}
\sqrt{\frac{1+x}{2}} X_{n-1}(x_k, x) &= \frac{(-1)^{k-1} \cos \left( n + \frac{1}{2} \right) \theta \sin \theta_k}{2(2n+1) \cos^3 \theta_{k/2} (\cos \theta - \cos \theta_k)} + \\
&+ \frac{(-1)^{k-1} \sin \left( n + \frac{1}{2} \right) \theta \sin \theta}{(2n+1) \sin \theta_k \cos \theta_{k/2} (\cos \theta - \cos \theta_k)} + \frac{(-1)^k 2n \cos \left( n + \frac{1}{2} \right) \theta}{(2n+1) \sin \theta_k \cos \theta_{k/2}} + \\
&+ \frac{(-1)^{k-1} \sin \left( n + \frac{1}{2} \right) \theta \sin \theta}{(2n+1) \sin \theta_k \cos \theta_{k/2} (\cos \theta - \cos \theta_k)} + \frac{(-1)^k \cos \left( n + \frac{1}{2} \right) \theta}{(2n+1) \sin \theta_k \cos \theta_{k/2}} + \\
&+ \frac{(-1)^{k-1} \cos \left( n + \frac{1}{2} \right) \theta}{2(2n+1) \sin \theta_k \cos \theta_{k/2}} \left[ \operatorname{cosec}^2 \frac{\theta + \theta_k}{2} + \operatorname{cosec}^2 \frac{\theta - \theta_k}{2} \right],
\end{aligned}$$

which proves the lemma.



LEMMA 4. We have writing  $x = \cos \theta$ ,  $x_k = \cos \theta_k$  etc. in (1.3)

$$\begin{aligned}
 \mu_k(x) = & \frac{\cos^4 \left( n + \frac{1}{2} \right) \theta \sin^4 \theta_k}{\left( n + \frac{1}{2} \right)^4 (\cos \theta - \cos \theta_k)^4} + \frac{(1 - 2 \cos \theta_k) \cos^4 \left( n + \frac{1}{2} \right) \theta \sin^2 \theta_k}{\left( n + \frac{1}{2} \right)^4 (\cos \theta - \cos \theta_k)^3} + \\
 & + \frac{\cos^4 \left( n + \frac{1}{2} \right) \theta \sin^4 \theta_k \cos \theta/2}{\left( n + \frac{1}{2} \right)^4 \cos^2 \theta_{k/2} (\cos \theta - \cos \theta_k)^3} + \\
 (2.6) \quad & + \frac{4 \cos^3 \left( n + \frac{1}{2} \right) \theta \sin \left( n + \frac{1}{2} \right) \theta \sin \theta \cos \theta/2 \sin^2 \theta_k}{\left( n + \frac{1}{2} \right)^3 (\cos \theta - \cos \theta_k)^3} + \\
 & + \frac{\cos^4 \left( n + \frac{1}{2} \right) \theta \sin^2 \theta_k \cos \theta/2}{\left( n + \frac{1}{2} \right)^4 (\cos \theta - \cos \theta_k)^2} \left[ \operatorname{cosec}^2 \frac{\theta + \theta_k}{2} + \operatorname{cosec}^2 \frac{\theta - \theta_k}{2} - \right. \\
 & \left. - \frac{4 \cos^4 \left( n + \frac{1}{2} \right) \theta \cos \theta/2 \sin^2 \theta_k}{\left( n + \frac{1}{2} \right)^3 (\cos \theta - \cos \theta_k)^2} \right].
 \end{aligned}$$

The proof of this lemma is straightforward and can be obtained by the preceding lemma.

LEMMA 5. If we break the summands in  $\mu_k(x)$  into two parts  $\mu_k^{(1)}(x)$  and  $\mu_k^{(2)}(x)$  in the following manner

$$\mu_k(x) = \mu_k^{(1)}(x) + \mu_k^{(2)}(x)$$

where,

$$\begin{aligned}
 \mu_k^{(1)}(x) = & \frac{4 \cos^3 \left( n + \frac{1}{2} \right) \theta \sin \left( n + \frac{1}{2} \right) \theta \sin \theta \cos \theta/2 \sin^2 \theta_k}{\left( n + \frac{1}{2} \right)^3 (\cos \theta - \cos \theta_k)^3} - \\
 & - \frac{4 \cos^4 \left( n + \frac{1}{2} \right) \theta \cos \theta/2 \sin^2 \theta_k}{\left( n + \frac{1}{2} \right)^3 (\cos \theta - \cos \theta_k)^2}
 \end{aligned}$$

and  $\mu_k^{(2)}(x)$  = the remaining ones, then we have

$$(2.7) \quad \sum_{k=1}^n |\mu_k^{(1)}(x)| \leq 120$$

and

$$(2.8) \quad \sum_{k=1}^n |(\cos \theta - \cos \theta_k) \mu_k^{(1)}(x)| \leq \frac{72 \cos \theta/2}{n}$$

$$(2.9) \quad \sum_{k=1}^n |\mu_k^{(2)}(x)| \leq 234$$

$$(2.10) \quad \sum_{k=1}^n \left| \sin \frac{1}{2} (\theta - \theta_k) \mu_k^{(2)}(x) \right| \leq \frac{130}{n}$$

and

$$(2.11) \quad \sum_{k=1}^n \sin^2 \frac{1}{2} (\theta - \theta_k) |\mu_k^{(2)}(x)| \leq \frac{78}{n^2}.$$

PROOF. (2.7) and (2.9) are straightforward if we recall a lemma proved by Saxena [4] whose version for our application is as follows. Let

$$E_k^m = \left[ \frac{\left| \cos \left( n + \frac{1}{2} \right) \theta \right|}{\left( n + \frac{1}{2} \right) \sin \frac{1}{2} |\theta - \theta_k|} \right]^m$$

then,

$$(2.12) \quad \sum_{k=1}^n E_k^m \leq 4 + 2^{m+1}, \quad \text{for } m = 2, 3, \dots$$

To establish (2.8), we see that

$$\begin{aligned} \sum_{k=1}^n |\cos \theta - \cos \theta_k| |\mu_k^{(1)}(x)| &\leq \sum_{k=1}^n \frac{4 \left| \cos^3 \left( n + \frac{1}{2} \right) \theta \sin \left( n + \frac{1}{2} \right) \theta \right| \sin \theta \sin^2 \theta_k}{\left( n + \frac{1}{2} \right)^3 (\cos \theta - \cos \theta_k)^2} + \\ &+ 4 \sum_{k=1}^n \frac{\cos^4 \left( n + \frac{1}{2} \right) \theta \cos \theta/2 \sin^2 \theta_k}{\left( n + \frac{1}{2} \right)^3 |\cos \theta - \cos \theta_k|} \leq \\ &\leq \frac{4 \cos \theta/2}{\left( n + \frac{1}{2} \right)} \sum_{k=1}^n \frac{\cos^2 \left( n + \frac{1}{2} \right) \theta}{\left( n + \frac{1}{2} \right)^3 \sin^2 \frac{1}{2} (\theta - \theta_k)} + \end{aligned}$$

$$\begin{aligned}
& + \frac{2 \cos \theta/2}{\left(n + \frac{1}{2}\right)} \sum_{k=1}^n \frac{\cos^4 \left(n + \frac{1}{2}\right) \theta}{\left(n + \frac{1}{2}\right)^2 \sin^2 \frac{1}{2} (\theta - \theta_k)} \equiv \\
& \equiv \frac{6 \cos \theta/2}{\left(n + \frac{1}{2}\right)} \sum_{k=1}^n E_k^2 < \frac{72 \cos \theta/2}{\left(n + \frac{1}{2}\right)}
\end{aligned}$$

which gives (2.8).

Now,

$$\begin{aligned}
\sum_{k=1}^n \sin^2 \frac{\theta - \theta_k}{2} |\mu_k^{(2)}(x)| & \equiv \sum_{k=1}^n \frac{\cos^4 \left(n + \frac{1}{2}\right) \theta}{\left(n + \frac{1}{2}\right)^4 \sin^2 \frac{1}{2} (\theta - \theta_k)} + \\
& + \frac{3}{2} \sum_{k=1}^n \frac{\cos^4 \left(n + \frac{1}{2}\right) \theta}{\left(n + \frac{1}{2}\right)^4 \sin^2 \frac{\theta - \theta_k}{2}} + \\
& + 2 \sum_{k=1}^n \frac{\cos^4 \left(n + \frac{1}{2}\right) \theta}{\left(n + \frac{1}{2}\right)^4 \sin^2 \frac{\theta - \theta_k}{2}} + 2 \sum_{k=1}^n \frac{\cos^4 \left(n + \frac{1}{2}\right) \theta}{\left(n + \frac{1}{2}\right)^4 \sin^2 \frac{\theta - \theta_k}{2}} \equiv \\
& \equiv \frac{13}{2 \left(n + \frac{1}{2}\right)^2} \sum_{k=1}^n E_k^2 < \frac{78}{n^2}
\end{aligned}$$

providing thereby (2.11).

LEMMA 6. For the polynomial  $\mu_k(x)$  in (1.3) we have for  $-1 \leq x \leq 1$ ,

$$(2.13) \quad \sum_{k=1}^n |x - x_k| |\mu_k'(x)| \leq 792.$$

PROOF. We can write

$$\begin{aligned}
& \sqrt{\frac{1+x}{2}} X_{n-1}(x_k, x) = \\
& = -\frac{1}{2(2n+1) \cos^3 \theta_{k/2}} \sum_{r=1}^{n-1} \cos \left(r + \frac{1}{2}\right) (\theta + \theta_k) + \cos \left(r + \frac{1}{2}\right) (\theta - \theta_k) - \\
& - \frac{2}{(2n+1) \sin \theta_k \cos \theta_{k/2}} \frac{d}{d\theta} \sum_{r=1}^{n-1} \cos \left(r + \frac{1}{2}\right) (\theta_k + \theta) + \cos \left(r + \frac{1}{2}\right) (\theta_k - \theta) =
\end{aligned}$$

$$\begin{aligned}
&= -\frac{1}{2(2n+1)\cos^3 \theta_{k/2}} \left[ \frac{\sin n(\theta+\theta_k)}{2 \sin \frac{\theta+\theta_k}{2}} + \frac{\sin n(\theta-\theta_k)}{2 \sin \frac{\theta-\theta_k}{2}} - \right. \\
(2.14) \quad &\left. -\cos \frac{1}{2}(\theta+\theta_k) - \cos \frac{1}{2}(\theta-\theta_k) \right] - \\
&\quad -\frac{2}{(2n+1)\sin \theta_k \cos \theta_{k/2}} \frac{d}{d\theta} \left[ \frac{\sin n(\theta+\theta_k)}{2 \sin \frac{\theta+\theta_k}{2}} + \frac{\sin n(\theta-\theta_k)}{2 \sin \frac{\theta-\theta_k}{2}} - \right. \\
&\quad \left. -\cos \frac{1}{2}(\theta+\theta_k) - \cos \frac{1}{2}(\theta-\theta_k) \right] = \\
&= -\frac{(-1)^{k-1} \cos \left( n + \frac{1}{2} \right) \theta \sin \theta_k}{2(2n+1)\cos^3 \theta \frac{k}{2} (\cos \theta_k - \cos \theta)} + \frac{\cos \theta/2 (-1)^{k-1}}{(2n+1)\cos^3 \theta_{k/2}} + \\
&\quad + \frac{2(-1)^{k-1}}{(2n+1)\cos \theta_{k/2}} \frac{d}{d\theta} \left[ \frac{\cos \left( n + \frac{1}{2} \right) \theta}{(\cos \theta_k - \cos \theta)} \right] + \frac{(-1)^{k-1} 4 \sin \theta/2}{(2n+1)\sin \theta_k}.
\end{aligned}$$

Owing to (2.14) we have by differentiating  $\mu_k(x)$  in Lemma 4,

$$\begin{aligned}
|\cos \theta - \cos \theta_k| |\mu'_k(\cos \theta)| &\leq \frac{122 \sin \left( n + \frac{1}{2} \right) \left( \frac{\theta - \theta_k}{2} \right)}{\left( n + \frac{1}{2} \right)^4 \sin^4 \frac{1}{2}(\theta - \theta_k)} + \\
&+ \frac{176 \left| \sin^3 \left( n + \frac{1}{2} \right) \left( \frac{\theta - \theta_k}{2} \right) \right|}{\left( n + \frac{1}{2} \right)^3 \sin^3 \frac{1}{2}|\theta - \theta_k|} + \frac{50 \sin^2 \frac{1}{2} \left( n + \frac{1}{2} \right) \left( \frac{\theta - \theta_k}{2} \right)}{\left( n + \frac{1}{2} \right)^2 \sin^2 \frac{1}{2}(\theta - \theta_k)} + \\
&+ \frac{6 \sin \left( n + \frac{1}{2} \right) \left( \frac{\theta - \theta_k}{2} \right)}{\left( n + \frac{1}{2} \right)^2 \sin^2 \frac{1}{2}|\theta - \theta_k|} + \frac{20}{n + \frac{1}{2}}
\end{aligned}$$

which explicitly gives the lemma.

**LEMMA 7.** For  $|x| \leq 1$

$$|R_n(f, 1) - R_n(f, x)| \leq 3046 \omega_f(1-x).$$

PROOF. From (1.8) we have

$$(2.15) \quad \begin{aligned} R_n(f, 1) - R_n(f, x) &= \sum_{k=1}^n [f(x_k) - f(1)] [\mu_k(1) - \mu_k(x)] + \\ &+ \frac{1-x}{2} [f(1) - f(-1)] \left[ 1 - \sum_{k=1}^n \mu_k(x) \right]. \end{aligned}$$

Hence, on account of Lemma 6 and the properties of modulus of continuity, we have,

$$\begin{aligned} |R_n(f, 1) - R_n(f, x)| &\leq 6\omega_f(1-x) + 708\omega_f(1-x) + 1416\omega_f(1-x) + \\ &+ 2\omega_f(1-x) \sum_{x_k \leq x} \left| (x - x_k) \frac{|\mu_k(x) - \mu_k(1)|}{(1-x)} \right| \leq \\ &\leq 2130\omega_f(1-x) + 2\omega_f(1-x) \sum_{k=1}^n |\xi - x_k| \mu'_k(\xi) \leq 3046\omega_f(1-x) \end{aligned}$$

which completes the proof.

### 3. Proof of the theorems

We have from (1.8) the identity

$$(3.1) \quad \begin{aligned} f(x) - R_n(f, x) &= \left[ \sum_{k=1}^n \mu_k(x) - 1 \right] [f(x) - f(x)] + \\ &+ \sum_{k=1}^n [f(x) - f(x_k)] \mu_k(x) = S_1 + S_2, \quad \text{say.} \end{aligned}$$

By virtue of

$$(3.2) \quad (1+x)\omega_f(1-x) + (1-x)\omega_f(1+x) \leq 6\omega_f(1-x^2), \quad x \in [-1, 1]$$

and Lemma 2, we obtain

$$(3.3) \quad |S_1| \leq 6\omega_f(1-x^2) \left[ \sum_{k=1}^n \mu_k(x) - 1 \right] \leq 36\omega_f \left( \frac{\sqrt{1-x^2}}{n} \right).$$

In order to estimate  $S_2$ , we break  $\mu_k(x)$  in the manner given in Lemma 5 and obtain the corresponding partitions of  $S_2$  as follows:

$$\begin{aligned} S_2^{(1)} &= \sum_{k=1}^n [f(x) - f(x_k)] \mu_k^{(1)}(x) \\ S_2^{(2)} &= \sum_{k=1}^n [f(x) - f(x_k)] \mu_k^{(2)}(x). \end{aligned}$$

For  $S_2^{(1)}$  we readily obtain,

$$(3.4) \quad \begin{aligned} |S_2^{(1)}| &\leq \omega_f \left( \frac{\cos \theta/2}{n} \right) \left[ \sum_{k=1}^n \left[ 1 + \frac{n(\cos \theta - \cos \theta_k)}{\cos \theta/2} |\mu^{(1)}(x)| \right] \right] \leq \\ &\leq 192\omega_f \left( \frac{\cos \theta/2}{n} \right), \quad \text{by Lemma 5.} \end{aligned}$$

To estimate  $S_2^{(2)}$  we make use the technique employed in [9] and notice that

$$(3.5) \quad |S_2^{(2)}| \leq 494\omega_f \left( \frac{\sqrt{1-x^2}}{n} + \frac{1}{n^2} \right).$$

Combination of (3.5)—(3.1) gives the proof.

Now to prove Theorem 2 we have to slightly modify the version of (2.11) which now reads as follows:

$$(3.6) \quad \sum_{k=1}^n \left| \sin^2 \frac{1}{2} (\theta - \theta_k) \mu_k^{(2)}(x) \right| \leq \frac{156 \cos \theta/2}{n},$$

which itself depends on the inequality

$$\left| \cos \left( n + \frac{1}{2} \right) \theta \right| \leq (2n+1) \cos \theta/2.$$

Finally, we sketch the proof of Theorem 3. We distinguish the two cases (i) when  $\sqrt{1-x^2} > \frac{1}{n}$  and (ii) when  $\sqrt{1-x^2} \leq \frac{1}{n}$  and consider them separately.

*Case 1.* If  $x \leq 0$ , then since

$$\sqrt{1+x} \leq \sqrt{1-x^2}$$

Theorem 2 directly implies Theorem 3.

Similarly we can consider the case  $x > 0$ .

*Case 2.* In this case, on account of (1.8) and (1.12), we have

$$\begin{aligned} f(x) - R_n^*(f, x) &= f(x) - R_n(f, x) - \frac{1+x}{2} [f(1) - R_n(f, 1)] = \\ &= \frac{1+x}{2} [f(x) - f(1)] + \frac{1+x}{2} [R_n(f, 1) - R_n(f, x)] + \frac{1-x}{2} [f(x) - R_n(f, x)]. \end{aligned}$$

Thus, owing to Lemma 6 and 7, we have,

$$\begin{aligned} |f(x) - R_n^*(f, x)| &\leq \\ &\leq \frac{1+x}{2} \omega_f(1-x) + 3046(1+x)\omega_f(1-x) + \frac{1-x}{2} \omega_f \left( \frac{\sqrt{1+x}}{n} \right) \leq \\ &\leq C_4 \omega_f(1-x^2) + C_5 \omega_f \left( \frac{1-x^2}{n} \right) \leq C_6 \omega_f \left( \frac{\sqrt{1-x^2}}{n} \right) \end{aligned}$$

providing thereby the proof of the theorem.

We can tacitly obtain the version of Theorem 3 for  $W_n(f, x)$  defined in (1.16) if we keep in mind Theorem 2 and the corresponding theorem for  $Q_n(f, x)$ .

ACKNOWLEDGEMENT. The author feels grateful to the referee for timely suggestions made by him.

## REFERENCES

- [1] FREUD, G., Egy Jackson-féle interpolációs eljárásról, *Mat. Lapok* **15** (1964), 330—336. *MR* 32#2794.
- [2] FREUD, G. and VÉRTESI, P., A new proof of A. F. Timan's theorem, *Studia Sci. Math. Hungar.* **2** (1967), 403—414. *MR* 36#4217.
- [3] MATHUR, K. K., On a proof of Jackson's theorem through interpolation process, *Studia Sci. Math. Hungar.* **6** (1971), 99—111. *MR* 49#5642.
- [4] SAXENA, R. B., On a polynomial of interpolation, *Studia Sci. Math. Hungar.* **2** (1967), 167—183. *MR* 35#4651.
- [5] SAXENA, R. B., The approximation of continuous functions by interpolatory polynomials, *Bulgar. Akad. Nauk. Otdel. Mat. Fiz. Nauk. Izv. Mat. Inst.* **12** (1970), 97—105. *MR* 43#7822.
- [6] SAXENA, R. B., A new proof of S. A. Telyakovskii's theorem on the approximation of continuous functions by algebraic polynomials, *Studia Sci. Math. Hungar.* **7** (1972), 3—9. *MR* 48#4575.
- [7] SAXENA, R. B., Approximation of continuous functions by polynomials, *Studia Sci. Math. Hungar.* **8** (1973), 437—446. *MR* 56#6197.
- [8] SALLAY, M., Über ein Interpolationsverfahren, *Magyar Tud. Akad. Mat. Kutató Int. Közl.* **9** (1964), 607—615. *MR* 32#4428.
- [9] SRIVASTAVA, K. B., A proof of Telyakovskii—Gopengauz's theorem through interpolation, *Serdica* **5** (1979), 272—279. *MR* 81e:41009.
- [10] TELYAKOVSKII, S. A., Two theorems on the approximation of functions by algebraic polynomials, *Math. Sbornik* **70** (1966), 252—255 (Russian). *MR* 33#1622.
- [11] VÉRTESI, P. and KIS, O., On a new interpolation process, *Ann. Univ. Sci. Budapest. Eötvös Sect. Math.* **10** (1967), 117—128. *MR* 40#4656.

(Received June 14, 1983)

DEPARTMENT OF MATHEMATICS  
INDIAN INSTITUTE OF TECHNOLOGY  
KANPUR—208016  
INDIA

## Present address

DEPARTMENT OF MATHEMATICS  
UNIVERSITY OF DAR ES SALAAM  
P.O. BOX 35062  
DAR ES SALAAM  
TANZANIA





## TESTS OF SUBHYPOTHESES IN LINEAR REGRESSION BASED ON RANK-ORDER ESTIMATES

CHING-YUAN CHIANG and MADAN L. PURI

### Summary

Under the linear regression model  $\underline{Y} = \alpha \underline{1}_n + \beta \underline{C} + \underline{Z}$ , a class of asymptotically distribution-free tests is proposed for testing the subhypothesis that some (but not all) components of the vector  $\beta$  of regression parameters are equal to 0. The tests are based on a normalized quadratic form in Jurečková's rank-order estimates of regression parameters, which is a natural analogue of the normal theory test statistic. The asymptotic efficiency of the proposed tests relative to the normal theory test is also examined.

### 1. Introduction

Consider the linear regression model

$$(1.1) \quad \underline{X}_n = \alpha \underline{1}_n + \beta \underline{C}_n + \underline{Z}_n$$

where

$$(1.2) \quad \underline{X}_n = (X_1, \dots, X_n)$$

is the random vector of observations,  $\alpha$  is an unknown scalar parameter (the intercept),

$$(1.3) \quad \underline{1}_n = (1, \dots, 1) \in \mathbf{R}^n,$$

$$(1.4) \quad \beta = (\beta_1, \dots, \beta_q)$$

is a  $q$ -dimensional vector of unknown regression parameters ( $q > 1$ ),

$$(1.5) \quad \underline{C}_n = (\underline{c}_1, \dots, \underline{c}_n)$$

is a  $q \times n$  matrix of known regression constants with columns

$$(1.6) \quad \underline{c}_i = (c_{1i}, \dots, c_{qi})' \quad (i = 1, \dots, n),$$

and

$$(1.7) \quad \underline{Z}_n = (Z_1, \dots, Z_n)$$

is the (unobservable) error random vector with the components independently

---

1980 *Mathematics Subject Classification*. Primary 62G10, 62G20; Secondary 62F10.

*Key words and phrases*. Linear regression model, subhypotheses, rank order estimates, asymptotic efficiencies.

distributed according to a common unknown distribution function

$$(1.8) \quad F(x) = P(Z_i \leq x) \quad (i = 1, \dots, n).$$

Let  $\beta$  be partitioned as

$$(1.9) \quad \beta = (\beta_1, \beta_2)$$

where

$$(1.10) \quad \beta_1 = (\beta_1, \dots, \beta_r), \quad \beta_2 = (\beta_{r+1}, \dots, \beta_q)$$

with  $1 \leq r < q$  being fixed, and consider the problem of testing the subhypotheses

$$(1.11) \quad H_0: \beta_2 = 0 \text{ vs. } H: \beta_2 \neq 0 ((\alpha, \beta_1) \text{ nuisance}).$$

Various versions of this problem, under the non-intercept form of the linear model, have been treated in detail in the classical normal theory (see, e.g., Graybill (1976), p. 194). Recently McKean and Hettmansperger (1976) have proposed a class of tests for (1.11) based on Jaeckel's (1972) dispersion measure, and Adichie (1978) has studied two class of aligned rank-order tests for a different version of (1.11), while Sen and Puri (1977) have proposed another class of aligned rank-order tests for a version of (1.11).

In the present paper we propose a class of asymptotically distribution-free tests for (1.11) based on a normalized quadratic form in Jurečková's (1971) rank-order estimates of regression parameters. The test statistics are natural analogues of the normal theory test statistic and have asymptotic chi-square distribution. The proposed tests are compared with the normal theory test for the same problem, and the asymptotic relative efficiency is derived.

## 2. Notations and assumptions

Let (1.1) be rewritten in the non-intercept form

$$(2.1) \quad X_n = \theta C_n^* + Z_n$$

where

$$(2.2) \quad \theta = (\alpha, \beta)$$

and

$$(2.3) \quad C_n^* = (1_n', C_n')'.$$

We make the usual assumption that the  $(q+1) \times n$  matrix  $C_n^*$  has the full rank  $q+1 < n$ . Let

$$(2.4) \quad \bar{c}_n = n^{-1} \sum_{i=1}^n c_i = (\bar{c}_{1n}, \dots, \bar{c}_{qn})'$$

where

$$(2.5) \quad \bar{c}_{mn} = n^{-1} \sum_{i=1}^n c_{mi} \quad (m = 1, \dots, q),$$

and define

$$(2.6) \quad D_n = \sum_{i=1}^n \zeta_i \zeta_i'.$$

Then the  $(q+1) \times (q+1)$  symmetric matrix

$$(2.7) \quad A_n = C_n^* C_n^{*'} = \begin{bmatrix} n & n\bar{\zeta}_n \\ n\bar{\zeta}_n & D_n \end{bmatrix}$$

has rank  $q+1$  and is positive definite. Consider the  $q \times q$  symmetric matrix

$$(2.8) \quad \begin{aligned} M_n &= \sum_{i=1}^n (\zeta_i - \bar{\zeta}_n)(\zeta_i - \bar{\zeta}_n)' = D_n - n\bar{\zeta}_n \bar{\zeta}_n' = \\ &= (\zeta_1 - \bar{\zeta}_n, \dots, \zeta_n - \bar{\zeta}_n)(\zeta_1 - \bar{\zeta}_n, \dots, \zeta_n - \bar{\zeta}_n)'. \end{aligned}$$

We note that by subtracting the first row of  $A_n$  left-multiplied by  $\bar{\zeta}_n$  from the second row block and then using (2.8) we obtain the matrix

$$(2.9) \quad \begin{bmatrix} n & n\bar{\zeta}_n' \\ 0' & M_n \end{bmatrix},$$

which, being row-equivalent to  $A_n$ , has rank  $q+1$  (see, e.g., Gantmacher (1959), Vol. 1, p. 45, Theorem 3). It follows that  $M_n$  has rank  $q$  and hence is positive definite. We also assume that the limiting matrices

$$(2.10) \quad A = \lim_{n \rightarrow \infty} n^{-1} A_n$$

and

$$(2.11) \quad M = \lim_{n \rightarrow \infty} n^{-1} M_n$$

exist and are positive definite.

Following Sen and Puri (1977), we simplify some of Jurečková's (1971) conditions on regression constants by assuming that each  $\zeta_i$  can be expressed as a difference

$$(2.12) \quad \zeta_i = \zeta_{i(1)} - \zeta_{i(2)}, \quad \zeta_{i(j)} = (c_{1i(j)}, \dots, c_{qi(j)})' \quad (i = 1, \dots, n; j = 1, 2)$$

where, for each  $m=1, \dots, q$  and each  $j=1, 2$ ,  $c_{mi(j)}$  is nondecreasing in  $i$ . We also assume that the  $c_{mi(j)}$  satisfy

$$(2.13) \quad \lim_{n \rightarrow \infty} n^{-1} \max_{1 \leq i \leq n} [c_{mi(j)} - \bar{c}_{mn(j)}]^2 = 0, \quad \bar{c}_{mn(j)} = n^{-1} \sum_{i=1}^n c_{mi(j)}$$

and

$$(2.14) \quad \lim_{n \rightarrow \infty} n^{-1} \sum_{i=1}^n [c_{mi(j)} - \bar{c}_{mn(j)}]^2 \in (0, \infty) \quad (m = 1, \dots, q; j = 1, 2).$$

Hence the  $c_{mi(j)}$  satisfy the Noether condition

$$(2.15) \quad \lim_{n \rightarrow \infty} \left\{ \max_{1 \leq i \leq n} [c_{mi(j)} - \bar{c}_{mn(j)}]^2 / \sum_{i=1}^n [c_{mi(j)} - \bar{c}_{mn(j)}]^2 \right\} = 0$$

$$(m = 1, \dots, q; j = 1, 2).$$

For  $\underline{b} = (b_1, \dots, b_q) \in \mathbf{R}^q$ , let

$$(2.16) \quad S_{nm}(\underline{b}) = \sum_{i=1}^n (c_{mi} - \bar{c}_{mn}) a_n[R_{ni}(\underline{b})], \quad (m = 1, \dots, q)$$

where

$$(2.17) \quad R_{ni} = \text{the rank of } X_i - \underline{b}\underline{c}_i \text{ among } X_1 - \underline{b}\underline{c}_1, \dots, X_n - \underline{b}\underline{c}_n \text{ in the ascending order } (i=1, \dots, n)$$

and the scores  $a_n(1), \dots, a_n(n)$  are generated by a non-constant, non-decreasing and square-integrable function  $\psi$  defined on  $(0, 1)$  in one of the following two ways:

$$(2.18) \quad a_n(i) = \psi(i/(n+1)), \quad (i = 1, \dots, n)$$

or

$$(2.19) \quad a_n(i) = E[\psi(U_{ni})], \quad (i = 1, \dots, n),$$

where  $U_{n1} \leq \dots \leq U_{nn}$  are the order statistics of a random sample of size  $n$  from the uniform distribution over  $(0, 1)$ . We further assume that  $\psi$  satisfies

$$(2.20) \quad \psi(1-u) = -\psi(u), \quad u \in (0, 1).$$

We note that such a function  $\psi$  satisfies

$$(2.21) \quad \int_0^1 \psi(u) du = 0$$

and

$$(2.22) \quad 0 < \lambda(\psi) = \left\{ \int_0^1 [\psi(u)]^2 du \right\}^{1/2} < \infty.$$

A class of score-generating functions of particular interest are of the form

$$(2.23) \quad \psi(u) = \varphi_g(u) = -g'[G^{-1}(u)]/g[G^{-1}(u)], \quad u \in (0, 1)$$

where  $G$  is a distribution function with a density  $g = G'$  which is symmetric and strongly unimodal and has a finite and positive Fisher information

$$(2.24) \quad 0 < I(g) = \int_{-\infty}^{\infty} [g'(x)/g(x)]^2 dG(x) < \infty.$$

Concerning the underlying distribution  $F$ , we assume that it has a symmetric and absolutely continuous density  $f = F'$  with finite Fisher information. We also assume that

$$(2.25) \quad \gamma(\psi, f) = \int_0^1 \psi(u) \varphi_f(u) du > 0,$$

where

$$(2.26) \quad \varphi_f(u) = -f'[F^{-1}(u)]/f[F^{-1}(u)], \quad u \in (0, 1).$$

### 3. The proposed tests

We first estimate  $\beta$ . Consider the set

$$(3.1) \quad B_n = \{b \in \mathbb{R}^q: \sum_{m=1}^q |S_{nm}(b)| = \text{minimum}\},$$

which is nonempty with probability one (see Jurečková (1971), Section 4). We choose one element

$$(3.2) \quad \hat{\beta}_n = (\hat{\beta}_{n1}, \dots, \hat{\beta}_{nq}) \in B_n$$

and partition it as

$$(3.3) \quad \hat{\beta}_n = (\hat{\beta}_{1n}, \hat{\beta}_{2n}), \quad \text{where} \quad \hat{\beta}_{2n} = (\hat{\beta}_{n,r+1}, \dots, \hat{\beta}_{nq}).$$

Let the scores  $a_n^+(1), \dots, a_n^+(n)$  be generated by the function

$$(3.4) \quad \psi^+(u) = \psi((u+1)/2), \quad u \in (0, 1)$$

according to

$$(3.5) \quad a_n^+(i) = \psi^+(i/(n+1)) \quad (i = 1, \dots, n)$$

or

$$(3.6) \quad a_n^+(i) = E[\psi^+(U_{ni})], \quad (i = 1, \dots, n).$$

For  $a \in \mathbb{R}$  let

$$(3.7) \quad W_n(a) = [n\lambda^2(\psi)]^{-1/2} \sum_{i=1}^n a_n^+[\hat{R}_{ni}(a)] \text{sign}(X_i - a - \hat{\beta}_{n1}c_i)$$

where

$$(3.8) \quad \begin{aligned} \hat{R}_{ni}(a) &= \text{the rank of } |X_i - a - \hat{\beta}_{n1}c_i| \text{ among} \\ &|X_1 - a - \hat{\beta}_{n1}c_1|, \dots, |X_n - a - \hat{\beta}_{n1}c_n|, \quad (i = 1, \dots, n) \end{aligned}$$

and  $\text{sign}(v)$  is equal to 1 or  $-1$  according as  $v \geq 0$  or  $v < 0$ . For arbitrarily fixed  $0 < \delta < 1$ , let  $z_{\delta/2}$  be the upper  $100(\delta/2)\%$  point of the standard normal distribution. Define

$$(3.9) \quad \alpha_n^* = \inf \{a \in \mathbb{R}: W_n(a) < z_{\delta/2}\}, \quad \alpha_n^{**} = \sup \{a \in \mathbb{R}: W_n(a) > -z_{\delta/2}\}.$$

Now let the matrix  $M_n$  (see (2.8)) be partitioned as

$$(3.10) \quad M_n = \begin{bmatrix} M_{n11} & M_{n12} \\ M_{n21} & M_{n22} \end{bmatrix}$$

where  $M_{n11}$  is  $r \times r$ , and define the  $(q-r) \times (q-r)$  matrix

$$(3.11) \quad \bar{M}_n = M_{n22} - M_{n21} M_{n11}^{-1} M_{n12},$$

which is symmetric and positive definite (because  $M_n$  is). Then a class of tests for (1.11) (indexed by the score-generating function  $\psi$ ) can be based on the normalized quadratic form

$$(3.12) \quad Q_n = [2Z_{\delta/2} n^{-1/2} (\alpha_n^{**} - \alpha_n^*)^{-1}]^2 \hat{\beta}_{2n} \bar{M}_n \hat{\beta}_{2n}'.$$

In anticipation of the comparison with the normal theory test (see Section 4), we consider a sequence of hypotheses

$$(3.13) \quad H_n: \beta_2 = n^{-1/2} b_2$$

where  $b_2 \in \mathbb{R}^{q-r}$  is arbitrary. We note that  $H_0$  is a special case of  $H_n$  where  $b_2 = 0$ . Let the matrix  $M$  (see (2.11)) be partitioned as

$$(3.14) \quad M = \begin{bmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{bmatrix}$$

where  $M_{11}$  is  $r \times r$ , and define the  $(q-r) \times (q-r)$  matrix

$$(3.15) \quad \bar{M} = M_{22} - M_{21} M_{11}^{-1} M_{12},$$

which is symmetric and positive definite. Then the asymptotic distribution of  $Q_n$  is as follows:

**THEOREM 1.** *Under  $H_n$ ,  $Q_n$  has asymptotically the non-central chi-square distribution  $\chi_{q-r}^2(\Delta_Q)$  with  $q-r$  degrees of freedom and the noncentrality parameter*

$$(3.16) \quad \Delta_Q = [\gamma(\psi, f)/\lambda(\psi)]^2 b_2 \bar{M} b_2'.$$

**COROLLARY 1.** *Under  $H_0$ ,  $Q_n$  has asymptotically the (central) chi-square distribution  $\chi_{q-r}^2$  with  $q-r$  degrees of freedom.*

For  $0 < \varepsilon < 1$  let  $\chi_{q-r, \varepsilon}^2$  be the upper  $100\varepsilon\%$  point of the  $\chi_{q-r}^2$  distribution. Then for large  $n$  we have the following asymptotically distribution-free test of approximately size  $\varepsilon$ :

$$(3.17) \quad \text{Reject } H_0 \text{ (in favour of } H) \text{ if and only if } Q_n \geq \chi_{q-r, \varepsilon}^2.$$

**PROOF OF THEOREM 1.** By (2.11), (3.10)–(3.11) and (3.14)–(3.15) we have

$$(3.18) \quad \lim_{n \rightarrow \infty} n^{-1} M_{njk} = M_{jk}; \quad j, k = 1, 2,$$

and

$$(3.19) \quad \lim_{n \rightarrow \infty} n^{-1} \bar{M}_n = \bar{M}.$$

By Theorem 4.1 of Jurečková (1971),  $n^{1/2}(\hat{\beta}_n - \beta) = (n^{1/2}(\hat{\beta}_{1n} - \beta_1), n^{1/2}(\hat{\beta}_{2n} - \beta_2))$  is asymptotically  $q$ -variate normal  $N_q(0, [\lambda(\psi)/\gamma(\psi, f)]^2 M^{-1})$ . So, by (3.14)–(3.15) and a well-known fact about inverses of partitioned matrices (see, e.g. Graybill (1976), p. 19, Theorem 1.3.1) we have

$$(3.20) \quad \mathcal{D}[n^{1/2}(\hat{\beta}_{2n} - \beta_2)] \rightarrow N_{q-r}(0, [\lambda(\psi)/\gamma(\psi, f)]^2 (\bar{M})^{-1})$$



where  $\mathcal{D}$  denotes distribution. In particular, under  $H_n$  we have

$$(3.21) \quad \mathcal{D}(n^{1/2} \hat{\beta}_{2n} | H_n) \rightarrow N_{q-r}(b_2, [\lambda(\psi)/\gamma(\psi, f)]^2 (\bar{M})^{-1}).$$

By Theorem 3.1 of McKean and Hettmansperger (1976),  $\gamma(\psi, f)$  is consistently estimated by

$$(3.22) \quad \gamma_n = 2z_{\delta/2} \hat{\lambda}(\psi) n^{-1/2} (\alpha_n^{**} - \alpha_n^*)^{-1}.$$

So, under  $H_n$ , the statistic

$$(3.23) \quad n^{1/2} \hat{\beta}_{2n} \gamma_n / \lambda(\psi) = 2z_{\delta/2} \hat{\beta}_{2n} (\alpha_n^{**} - \alpha_n^*)^{-1}$$

is asymptotically  $N_{q-r}(\gamma(\psi, f) b_2 / \lambda(\psi), (\bar{M})^{-1})$ . By (3.12) and (3.23) we can express  $Q_n$  as

$$(3.24) \quad Q_n = [n^{1/2} \hat{\beta}_{2n} \gamma_n / \lambda(\psi)] n^{-1} \bar{M}_n [n^{1/2} \hat{\beta}_{2n} \gamma_n / \lambda(\psi)]'.$$

So, by (3.19),  $Q_n$  under  $H_n$  is asymptotically non-central chi-square with  $q-r$  degrees of freedom and noncentrality parameter

$$[\gamma(\psi, f) b_2 / \lambda(\psi)] \bar{M} [\gamma(\psi, f) b_2 / \lambda(\psi)]' = \Delta_Q.$$

REMARKS. 1. By (3.24), the test statistic can be expressed as

$$(3.25) \quad Q_n = (\hat{\beta}_{2n} \bar{M}_n \hat{\beta}_{2n}' / [\lambda(\psi) / \gamma_n]^2$$

where  $[\lambda(\psi) / \gamma_n]^2$  is a consistent estimate of the coefficient  $[\lambda(\psi) / \gamma(\psi, f)]^2$  of the asymptotic covariance matrix of  $n^{1/2} \hat{\beta}_{2n}$ . Thus  $Q_n$  is a natural analogue of the normal theory test statistic (see (4.11) and (4.30)).

2. Among the subclass of score-generating functions given by (2.23), if  $G$  is the standard logistic distribution function, then  $\psi(u) = \varphi_\theta(u) = 2u - 1$  and  $\psi^+(u) = u$ , which generate Wilcoxon-type scores; and if  $G = \Phi$  is the standard normal distribution function, then  $\psi(u) = \varphi_\theta(u) = \Phi^{-1}(u)$  and  $\psi^+(u) = \Phi^{-1}[(u+1)/2]$ , which generate normal scores.

#### 4. Asymptotic efficiency

For the purpose of comparison with the classical normal theory test of (1.11), we make the additional assumption that  $F$  has a finite and positive (but unknown) variance

$$(4.1) \quad 0 < \text{Var}(Z_i) = \sigma^2 < \infty \quad (i = 1, \dots, n).$$

We note that

$$(4.2) \quad E(Z_i) = 0, \quad (i = 1, \dots, n).$$

Consider the least-squares estimate of  $\theta$  based on  $X_n$

$$(4.3) \quad \tilde{\theta}_n = (\tilde{z}_n, \tilde{\beta}_n) = X_n C_n^{*'} A_n^{-1}$$

and the corresponding unbiased estimate of  $\sigma^2$

$$(4.4) \quad s_n^2 = (\underline{X}_n - \underline{\tilde{\theta}}_n C_n^*)(\underline{X}_n - \underline{\tilde{\theta}}_n C_n^*)' / (n - q - 1).$$

Let  $\underline{\tilde{\beta}}_n = (\underline{\tilde{\beta}}_{n1}, \dots, \underline{\tilde{\beta}}_{nq})$  be partitioned as

$$(4.5) \quad \underline{\tilde{\beta}}_n = (\underline{\tilde{\beta}}_{1n}, \underline{\tilde{\beta}}_{2n}), \quad \text{where} \quad \underline{\tilde{\beta}}_{2n} = (\underline{\tilde{\beta}}_{n, r+1}, \dots, \underline{\tilde{\beta}}_{nq}).$$

We also partition  $A_n$  as

$$(4.6) \quad A_n = \begin{bmatrix} A_{n11} & A_{n12} \\ A_{n21} & A_{n22} \end{bmatrix}$$

where  $A_{n11}$  is  $(r+1) \times (r+1)$ , and define

$$(4.7) \quad \bar{A}_n = A_{n22} - A_{n21} A_{n11}^{-1} A_{n12}.$$

Then, the normal theory test of (1.11) is based on the  $F$ -statistic

$$(4.8) \quad \mathcal{F}_n = \underline{\tilde{\beta}}_{2n}' \bar{A}_n \underline{\tilde{\beta}}_{2n} / (q - r) s_n^2$$

(see, e.g., Anderson (1971), Section 2.2), or equivalently on the statistic

$$(4.9) \quad L_n = (q - r) \mathcal{F}_n = \underline{\tilde{\beta}}_{2n}' \bar{A}_n \underline{\tilde{\beta}}_{2n} / s_n^2.$$

It is well-known that if  $F$  is normal, then  $\mathcal{F}_n$  under  $H_0$  has the  $F$ -distribution with  $q - r$  and  $n - q - 1$  degrees of freedom. It will be shown later (in the proof of Theorem 2) that

$$(4.10) \quad \bar{A}_n = \bar{M}_n.$$

Thus  $L_n$  can also be expressed as

$$(4.11) \quad L_n = \underline{\tilde{\beta}}_{2n}' \bar{M}_n \underline{\tilde{\beta}}_{2n} / s_n^2.$$

The following theorem gives the asymptotic distribution of  $L_n$  under  $H_n$  but under no assumption concerning the specific shape of  $F$ .

**THEOREM 2.** *Under  $H_n$ ,  $L_n$  is asymptotically  $\chi_{q-r}^2(\Delta_L)$ , where*

$$(4.12) \quad \Delta_L = \sigma^{-2} b_2 \bar{M} b_2'.$$

To compare the proposed  $Q_n$ -tests with the normal theory test, we make the further assumption that

$$(4.13) \quad b_2 \neq 0,$$

which, by the positive definiteness of  $\bar{M}$ , makes the right-hand sides of (4.12) and (3.16) strictly positive. Combining Theorems 1 and 2, we have the following result regarding the asymptotic relative efficiency.

COROLLARY 2. *The asymptotic relative efficiency of the  $Q_n$ -tests with respect to the normal theory test (based on  $L_n$ ) of (1.11) is*

$$\begin{aligned}
 e_{Q,L}(F) &= \sigma^2 [\gamma(\psi, f) / \lambda(\psi)]^2 \\
 (4.14) \quad &= \sigma^2 \left[ \int_0^1 \psi(u) \phi_f(u) du \right]^2 / \int_0^1 [\psi(u)]^2 du \\
 &= \sigma^2 \left[ \int_0^1 \psi^+(u) \phi_f((u+1)/2) du \right]^2 / \int_0^1 [\psi^+(u)]^2 du.
 \end{aligned}$$

The quantity given in (4.14) has been extensively studied. For example, if  $\psi(u) = 2u - 1$  (i.e.,  $\psi^+(u) = u$ ), then  $e_{Q,L}(F)$  has a lower bound of 0.864, and is equal to  $3/\pi$  ( $\doteq 0.955$ ) if the underlying distribution  $F$  is normal. And if  $\psi = \Phi^{-1}$ , then  $e_{Q,L}(F)$  is not less than 1, and is strictly greater than 1 unless  $F$  is normal (see, e.g., Puri and Sen (1971), Section 3.8).

PROOF OF THEOREM 2. We first establish (4.10). Let each  $\underline{\zeta}_i$  (see (1.6)) be partitioned as

$$(4.15) \quad \underline{\zeta}_i = (\underline{\zeta}'_{i,1}, \underline{\zeta}'_{i,2})'$$

where

$$(4.16) \quad \underline{\zeta}_{i,1} = (c_{1i}, \dots, c_{ri})', \quad \underline{\zeta}_{i,2} = (c_{r+1,i}, \dots, c_{qi})',$$

and let

$$(4.17) \quad \underline{\zeta}_{n,j} = n^{-1} \sum_{i=1}^n \underline{\zeta}_{i,j} \quad (j = 1, 2).$$

Then, by (2.4)–(2.7) and (4.6) we have

$$(4.18) \quad A_{n11} = \begin{bmatrix} n & n\bar{\zeta}'_{n,1} \\ n\bar{\zeta}_{n,1} & \sum_{i=1}^n \underline{\zeta}_{i,1} \underline{\zeta}'_{i,1} \end{bmatrix}$$

$$(4.19) \quad A_{n22} = \sum_{i=1}^n \underline{\zeta}_{i,2} \underline{\zeta}'_{i,2},$$

$$(4.20) \quad A_{n21} = \left[ \sum_{i=1}^n \underline{\zeta}_{i,2}, \sum_{i=1}^n \underline{\zeta}_{i,2} \underline{\zeta}'_{i,1} \right]$$

and

$$(4.21) \quad A_{n12} = A'_{n21}.$$

On the other hand, by (2.4)–(2.5), (2.8), (3.10) and (4.15)–(4.17) we have

$$(4.22) \quad M_{njk} = \sum_{i=1}^n (\underline{\zeta}_{i,j} - \bar{\zeta}_{n,j})(\underline{\zeta}_{i,k} - \bar{\zeta}_{n,k})' = \sum_{i=1}^n \underline{\zeta}_{i,j} \underline{\zeta}'_{i,k} - n\bar{\zeta}_{n,j} \bar{\zeta}'_{n,k} \quad (j, k = 1, 2).$$

By direct computation and (3.22) (with  $j=k=1$ ), it is easily checked that inverse of  $A_{n11}$  is

$$(4.23) \quad A_{n11}^{-1} = \begin{bmatrix} \frac{1}{n} + \bar{c}'_{n,1} M_{n11}^{-1} \bar{c}_{n,1} & -\bar{c}'_{n,1} M_{n11}^{-1} \\ -M_{n11}^{-1} \bar{c}_{n,1} & M_{n11}^{-1} \end{bmatrix}.$$

By further routine computation and (4.22), we have

$$(4.24) \quad A_{n21} A_{n11}^{-1} A_{n12} = n \bar{c}_{n,2} \bar{c}'_{n,2} + M_{n21} M_{n11}^{-1} M_{n12},$$

from which by (3.11), (4.7), (4.19) and (4.22) (with  $j=k=2$ ) we obtain (4.10).

Let the matrix  $A$  (see (2.10)) be partitioned as

$$(4.25) \quad A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

where  $A_{11}$  is  $(r+1) \times (r+1)$ , and define

$$(4.26) \quad \bar{A} = A_{22} - A_{21} A_{11}^{-1} A_{12}.$$

Then, by (4.6)–(4.7) we have

$$(4.27) \quad \lim_{n \rightarrow \infty} n^{-1} A_{njk} = A_{jk} \quad (j, k = 1, 2)$$

and

$$(4.28) \quad \lim_{n \rightarrow \infty} n^{-1} \bar{A}_n = \bar{A}.$$

It follows from (3.19) and (4.10) that

$$(4.29) \quad \bar{A} = \bar{M}.$$

Now, under our regularity assumptions, the random vector

$$n^{1/2} \bar{\theta}_n = n^{1/2} [(\bar{\alpha}_n, \bar{\beta}_n) - (\alpha, \beta)] = (n^{1/2} [(\bar{\alpha}_n, \bar{\beta}_{1n}) - (\alpha, \beta_1)], n^{1/2} (\bar{\beta}_{2n} - \beta_2))$$

is asymptotically  $(q+1)$ -variate normal  $N_{q+1}(0, \sigma^2 A^{-1})$ , and  $s_n^2$  is a consistent estimate of  $\sigma^2$  (see Anderson (1971), p. 25, Corollary 2.6.1 and Theorem 2.6.2). Consequently, by (4.25)–(4.26) we have

$$\mathcal{D}[n^{1/2} (\bar{\beta}_{2n} - \beta_2)] \rightarrow N_{q-r}(0, \sigma^2 (\bar{A})^{-1})$$

and hence

$$(4.30) \quad \mathcal{D}(n^{1/2} \bar{\beta}_{2n} | H_n) \rightarrow N_{q-r}(b_2, \sigma^2 (\bar{A})^{-1}).$$

Thus, by (4.9), (4.28)–(4.29) and with  $s_n^2$  as a consistent estimator of  $\sigma^2$ ,  $L_n$  under  $H_n$  is asymptotically non-central chi-square with  $q-r$  degrees of freedom and noncentrality parameter  $\sigma^{-2} b_2' \bar{A} b_2 = \sigma^{-2} b_2' \bar{M} b_2 = \Delta_L$ . The proof is completed.

## REFERENCES

- [1] ADICHIE, J. N., Rank tests of sub-hypotheses in the general linear regression, *Ann. Statist.* 6 (1978), 1012—1026. *MR* 80a: 62053.
- [2] ANDERSON, T. W., *The Statistical Analysis of Time Series*, J. Wiley, New York, 1971. *MR* 44#1169.
- [3] GANTMACHER, F. R., *The Theory of Matrices*, Chelsea, New York, 1959. *MR* 21#6372c.
- [4] GRAYBILL, F. A., *Theory and Application of the Linear Model*, Duxbury Press, North Scituate, Mass. 1976. *MR* 56#13457.
- [5] JAECKEL, L. A., Estimating regression coefficients by minimizing the dispersion of the residuals, *Ann. Math. Statist.* 43 (1972), 1449—1458. *MR* 50#1424.
- [6] JUREČKOVÁ, J., Nonparametric estimate of regression coefficients, *Ann. Math. Statist.* 42 (1971), 1328—1338. *MR* 45#4553.
- [7] MCKEAN, J. W. and HETTMANSPERGER, T. P., Tests of hypotheses based on ranks in the general linear model, *Comm. Statist. — Theory Methods* A5 (1976), 693—709. *MR* 55#9399.
- [8] PURI, M. L. and SEN, P. K., *Nonparametric Methods in Multivariate Analysis*, J. Wiley, New York—London—Sidney, 1971. *MR* 45#7893.
- [9] SEN, P. K. and PURI, M. L., Asymptotically distribution-free aligned rank order tests for composite hypotheses for general multivariate linear models, *Z. Wahrscheinlichkeitstheorie und Verw. Gebiete* 39 (1977), 175—186. *MR* 57#17936.

(Received June 17, 1983)

DEPARTMENT OF MATHEMATICS AND  
COMPUTER SCIENCE  
JAMES MADISON UNIVERSITY  
HARRISONBURG, VA 22807

DEPARTMENT OF MATHEMATICS  
SWAIN HALL EAST  
INDIANA UNIVERSITY  
BLOOMINGTON, IN 47405  
U.S.A.



# REMARKS ON COMBINATORIAL GEOMETRY I

JÓZSEF BECK

## 1. Introduction

In this note we prove two results in “geometric configuration calculus”. First let there be given a set  $S$  of  $n$  distinct points in the Euclidean plane. We say that the straight line  $L$  is determined by  $S$  if  $L$  contains at least two points of  $S$ , i.e.,  $|L \cap S| \geq 2$  (throughout this paper  $|H|$  denotes, as usual, the number of elements of the set  $H$ ). Let  $\text{line}(S)$  denote the number of straight lines determined by  $S$ .

The complete determination of the set of pairs  $(|S|, \text{line}(S))$  seems to be hopeless at present, but some interesting results are known about it. It is easily seen that the values  $\text{line}(S) = \binom{n}{2} - 1$  and  $\binom{n}{2} - 3$  are impossible for  $|S| = n$ . Answering a question of B. Grünbaum, P. Erdős [2] proved the following “jumping point” phenomenon in the distribution of  $\text{line}(S)$ .

**THEOREM A** (P. Erdős). *There are positive absolute constants  $c_1, c_2$  and  $c_3$  such that, if the integer  $l$  belongs to the interval  $\left[c_1 n^{3/2}, \binom{n}{2} - 4\right]$  then there is an  $n$ -element point-set  $S$  with  $\text{line}(S) = l$ . On the other hand, there exists a “gap”  $[N, N + c_2 n]$  where  $N > c_3 n^{3/2}$  such that  $\text{line}(S) \notin [N, N + c_2 n]$  whenever  $|S| = n$ .*

At the end of his paper Erdős raised the problem of finding analogous result for circles.

We say that the circle  $R$  is determined by the point-set  $S$  if  $R$  contains at least three points of  $S$ . Denote by  $\text{circle}(S)$  the number of circles determined by  $S$ .

Our first object is to answer Erdős’ question as follows.

**THEOREM 1.1.** *There are positive absolute constants  $c_4, c_5$  and  $c_6$  such that, if the integer  $k$  belongs to the interval  $\left[c_4 n^{5/2}, \binom{n}{3} - 38\right]$  then there is an  $n$ -element point-set  $S$  with  $\text{circle}(S) = k$ . On the other hand, there is a “gap”  $[N, N + c_5 n^2]$  where  $N > c_6 n^{5/2}$  such that  $\text{circle}(S) \notin [N, N + c_5 n^2]$  whenever  $|S| = n$ .*

Roughly speaking our theorem states that the “jumping point” in the distribution of  $\text{circle}(S)$  is about  $n^{5/2}$  ( $|S| = n$ ).

The reader is encouraged to consult T. Zaslavsky’s survey paper [6] of recent progress in related problems and results.



Second, let there be given a set  $S = \{P_1, P_2, \dots, P_n\}$  of  $n$  distinct points in general position (i.e., no three of them are on the same straight line) in the Euclidean plane, and consider the  $\binom{n}{2}$  straight-line segments  $\overline{P_i P_j}$  with endpoints  $P_i, P_j$  ( $1 \leq i < j \leq n$ ). Let  $\text{cross}(S)$  denote the number of points (counted with multiplicity) at which the segments  $\overline{P_i P_j}$  ( $1 \leq i < j \leq n$ ) cross each other. Let  $\text{cross}^*(S)$  denote the number of *distinct* points at which the segments  $\overline{P_i P_j}$  ( $1 \leq i < j \leq n$ ) cross each other.

It is well-known that  $\text{cross}(S) = \binom{n}{4}$  whenever  $S$  forms a convex polygon of  $n$  vertices. Indeed, then any four points of  $S$  determines exactly one cross-point. Several years ago P. Erdős (oral communication) suspected that  $\text{cross}^*(S)$  cannot be much less than  $\binom{|S|}{4}$ . Here we shall prove the following result.

**THEOREM 1.2.** *For any set  $S$  of  $n$  distinct points in general position in the plane,*

$$\text{cross}^*(S) > c_7 \binom{n}{4}$$

where  $c_7 > 0$  is independent of  $n$ .

For a survey concerning *crossing numbers* of graphs see e.g. Grünbaum [3].

## 2. Proof of Theorem 1.1

We start with the constructive part. Let there be given two concentric circles  $R_1$  and  $R_2$  with common center  $C$ , and points  $P_i \in R_1$ ,  $-m \leq i \leq m$ ,  $Q_j \in R_2$ ,  $-M \leq j \leq M$  (here  $m \leq M$ ) with the properties

- ( $\alpha$ ) the points  $P_{-m}, P_{-m+1}, \dots, P_0, P_1, \dots, P_m$  are equidistant on  $R_1$  in this order, and similarly,  $Q_{-M}, Q_{-M+1}, \dots, Q_0, Q_1, \dots, Q_M$  are equidistant on  $R_2$  in this order;
- ( $\beta$ ) the half-line starting from the common center  $C$  and passing through  $P_i$  contains  $Q_i$ ,  $-m \leq i \leq m$  ( $m$  and  $M$  will be specified later depending on the value of  $k$ ).

Denote by  $f(2M+1, 2m+1)$  the number of circles  $R$  such that

$$|R \cap \{P_i : -m \leq i \leq m\}| = 2 \text{ and } |R \cap \{Q_j : -M \leq j \leq M\}| = 2.$$

An easy calculation gives

$$(1) \quad f(2M+1, 2m+1) = \binom{2m+1}{2} M - 4 \binom{m+1}{3}.$$

We introduce the concept of Defect ( $S$ ) as follows

$$D(S) = \text{Defect}(S) = \binom{|S|}{3} - \text{circle}(S).$$

Set

$$S_1 = S_1(m) = \{P_i: -m \leq i \leq m\} \quad \text{and} \quad S_2 = S_2(M) = \{Q_j: -M \leq j \leq M\}.$$

Clearly

$$(2) \quad D(S_1 \cup S_2) = \binom{2M+1}{3} - 1 + \binom{2m+1}{3} - 1 + 3 \cdot f(2M+1, 2m+1).$$

If  $c_4$  is sufficiently large and  $k \in \left[ c_4 n^{5/2}, \binom{n}{3} - 38 \right]$  then from (1) and (2) it follows by elementary calculation the existence of integers  $M$  and  $m$  ( $M \geq m \geq 0$ ) such that

$$(3) \quad 38 \leq \binom{n}{3} - k - D(S_1(m) \cup S_2(M)) \leq \frac{1}{2} \binom{y}{3}$$

where  $y = n - (2M+1) - (2m+1)$  and  $y \geq n^{1/2}$ . Indeed, let  $M$  be the largest integer with  $\binom{2M+1}{3} - 1 \leq \binom{n}{3} - 38 - k$ , and next let  $m$  be the largest integer with  $\binom{2m+1}{3} - 1 + 3f(2M+1, 2m+1) \leq \binom{n}{3} - 38 - k - \binom{2M+1}{3} + 1$ .

We require the following simple lemma.

LEMMA 2.1. If  $k^* \in \left[ \frac{1}{2} \binom{y}{3}, \binom{y}{3} - 38 \right]$  and  $y$  is sufficiently large then there exists a point-set  $Y$  such that  $|Y| = y$  and  $\text{circle}(Y) = k^*$ .

The proof of Lemma 2.1 proceeds along the same lines as Erdős argued in his original paper. For the sake of completeness we include it. Similarly as above, one can easily find nonnegative integers  $u, v, z$  in this order such that

$$y - (u+v+z) > \frac{1}{6} y \quad \text{and} \quad 38 \leq \binom{y}{3} - k^* - \left\{ \binom{u}{3} - 1 + \binom{v}{3} - 1 + \binom{z}{3} - 1 \right\} = O(y^{8/9}).$$

Moreover, for any integer  $w$  with  $38 \leq w = O(y^{8/9})$  there must exist nonnegative integers  $p, q$  and  $r$  such that  $p+q+r < \frac{1}{6} y$  and  $3p+9q+19r=w$ .

Now the construction of the desired  $y$ -element set  $Y$  goes as follows. Consider three sets  $U, V$  and  $Z$  having the properties  $|U|=u, |V|=v, |Z|=z, U \subset R^{(1)}, V \subset R^{(2)}, Z \subset R^{(3)}$  where  $R^{(i)}, i=1, 2, 3$  are different circles. Let  $H_4^{(i)}, 1 \leq i \leq p$  be four-element sets such that  $H_4^{(i)} \subset R_4^{(i)}$  where the circles  $R_4^{(i)}, i=1, 2, \dots, p$  are different. Let  $H_5^{(j)}, 1 \leq j \leq q$  and  $H_6^{(l)}, 1 \leq l \leq r$  be five-element and six-element sets, respectively, defined analogously. Now let

$$Y = U \cup V \cup Z \cup \bigcup_{i=1}^p H_4^{(i)} \cup \bigcup_{j=1}^q H_5^{(j)} \cup \bigcup_{l=1}^r H_6^{(l)} \cup A$$

where  $|A| = y - (u+v+z+p+q+r)$ .

In the construction above one can easily guarantee the equality below

$$\begin{aligned} D(Y) &= D(U) + D(V) + D(Z) + \sum_{i=1}^p D(H_4^{(i)}) + \sum_{j=1}^q D(H_5^{(j)}) + \sum_{l=1}^r D(H_6^{(l)}) = \\ &= \binom{u}{3} - 1 + \binom{v}{3} - 1 + \binom{z}{3} - 1 + 3p + 9q + 19r = \binom{y}{3} - k^*, \end{aligned}$$

which completes the proof of Lemma 2.1.  $\square$

By Lemma 2.1 and (3) there exist a  $y$ -element set  $Y$  such that  $D(Y) = \binom{n}{3} - k - D(S_1 \cup S_2)$ . Observe that if  $y$  is small then we are immediately done. Indeed, then  $n$  is also small (since  $y \geq n^{1/2}$ ) and so the interval  $\left[ c_4 n^{5/2}, \binom{n}{3} - 38 \right]$  is empty.

Moreover, we may assume that  $Y$  is in general position relative to  $S_1 \cup S_2$ , i.e.,

$$D(S_1 \cup S_2 \cup Y) = D(S_1 \cup S_2) + D(Y).$$

Summarizing, we have

$$k = \binom{n}{3} - D(S_1 \cup S_2 \cup Y),$$

that is, the  $n$ -element point-set  $S = S_1 \cup S_2 \cup Y$  determines exactly  $k$  circles.

To prove the existence of a large "gap"  $[N, N + c_5 n^2]$  with some  $N > c_6 n^{5/2}$  we need the following very recent result conjectured by Erdős and proved independently of each other by Szemerédi and Trotter [5] and the author [1].

**THEOREM B.** *Let there be given  $n$  points in the plane, no  $n-t$  on the same straight line ( $t$ ,  $0 \leq t \leq n-3$  is arbitrary). Then the number of straight lines containing at least two of the given points exceeds  $c_0 t n$  ( $c_0 > 0$  independent of  $n$  and  $t$ ).*

Applying inversions with centers at the given points one can easily deduce from Theorem B the analogous result for circles.

**LEMMA 2.2.** *Let there be given  $n$  points in the plane, no  $n-t$  on the same circle (straight line). Then the number of circles (straight lines) containing at least three of the points exceeds  $c_0^* t n^2$ .  $\square$*

Now let there be given an  $n$ -element set  $S$  in the plane. Assume that some circle (straight line) contains exactly  $n-x$  points of  $S$ . It is easily seen that circle  $(S) \in I(x)$  where

$$I(x) = \left[ 1 + \binom{n-x}{2} x - (n-x) \binom{x}{2}, 1 + \binom{n-x}{2} x + (n-x) \binom{x}{2} + \binom{x}{3} \right].$$

By some elementary calculation it follows that there exists an interval  $[N, N + c_5 n^2]$  with  $N = c n^{5/2}$  such that the union  $\bigcup_{1 \leq x \leq n/10} I(x)$  of the first  $\frac{n}{10}$  intervals  $I(x)$ ,  $1 \leq x \leq \frac{n}{10}$  is disjoint from  $[N, N + c_5 n^2]$ . This means that there is no  $S$  such that

$|S|=n$ , some circle (straight line) contains  $\cong \frac{9}{10}n$  points of  $S$  and  $\text{circle}(S) \in [N, N+c_5n^2]$ .

Now we are ready, since in the opposite case (i.e., if any circle or straight line contains  $< \frac{9}{10}n$  of the points) Lemma 2.2 immediately yields  $\text{circle}(S) > cn^3$ . Theorem 1.1 follows.  $\square$

We remark that instead of Lemma 2.2 it suffices to use some old results of Kelly and Moser [4]. They proved that if  $n$  points are given in the plane, with no  $n-t$  points collinear, and if  $n > c_8 t^2$  then the number  $l$  of straight lines containing two or more of the points satisfies  $l \geq tn - c_9 t^3$ . Kelly and Moser also show that if  $l_i$  is the number of straight lines containing exactly  $i$  of the points, then

$$(4) \quad l_2 \geq 3 + \sum_{i \geq 4} (i-3)l_i.$$

Adding  $l_2 + l_3$  to both sides of (4) we obtain

$$2(l_2 + l_3) \geq 2l_2 + l_3 \geq 3 + l_2 + l_3 + l_4 + \dots = 3 + l.$$

Consequently

$$l_2 + l_3 \geq 3/2 + l/2.$$

Summarizing, if  $n$  points are given in the plane, with no  $n-t$  points collinear, and if  $n > c_8 t^2$  then the number of straight lines  $l_2 + l_3$  containing two or three points satisfies  $l_2 + l_3 > tn/2 - c_9 t^3/2$ .

Applying now inversions we get a weaker version of Lemma 2.2 working for  $n > c_8 t^2$  only. But it is sufficiently strong to complete the proof of Theorem 1.1.

### 3. Proof of Theorem 1.2

Let there be given  $n$  points  $P_1, P_2, \dots, P_n$  in general position in the plane. Let  $\overline{P_i P_j}$  ( $1 \leq i < j \leq n$ ) denote the straight-line segments with endpoints  $P_i, P_j$ , and denote by  $L(P_i, P_j)$  the straight line containing  $P_i$  and  $P_j$  (clearly  $\overline{P_i P_j} \subset L(P_i, P_j)$ ). Let  $Q_1, Q_2, \dots, Q_r$  be the points at which the segments  $\overline{P_i P_j}$  ( $1 \leq i < j \leq n$ ) cross each other (i.e.,  $r = \text{cross}^*(\{P_1, \dots, P_n\})$ ). Moreover, let  $m_k$  denote the number of segments  $\overline{P_i P_j}$  passing through the point  $Q_k$ ,  $1 \leq k \leq r$  (multiplicity). Since any five points of  $P_1, \dots, P_n$  determines at least one cross-point  $Q_k$ , we get

$$(5) \quad \sum_{k=1}^r \binom{m_k}{2} \geq \binom{n}{5} / (n-4).$$

Furthermore, observe that

$$(6) \quad m_k \leq n/2, \quad 1 \leq k \leq r.$$

Now we recall the dual form of Theorem 1.5 in [1].

LEMMA 3.1. Let there be given  $N$  straight lines  $L_1, \dots, L_N$  in the plane, and let  $Q_1, \dots, Q_R$  denote their intersection points. Let  $M_k$  denote the number of straight lines  $L_i$  passing through  $Q_k$ ,  $1 \leq k \leq R$  (multiplicity). Then for every positive  $M \leq \sqrt{2N}$ ,  $\sum_{k: M_k \geq M} 1 < c_{10} \frac{N^2}{M^{2+\delta}}$  with some positive absolute constant  $\delta$ .  $\square$

Note that in [1] we proved  $\delta \geq 1/20$ , but Szemerédi and Trotter [5] succeeded in determining the exact value of  $\delta$ , namely  $\delta = 1$ .

Applying Lemma 3.1 to the lines  $L(P_i, P_j)$ ,  $1 \leq i < j \leq n$  (and so  $N = \binom{n}{2}$ ), we obtain (we also use (6))

$$(7) \quad \sum_{k: m_k \geq M} 1 \leq \sum_{i=0}^{\infty} \sum_{2^i M \leq m_k < 2^{i+1} M} c_{10} \frac{\binom{n}{2}^2}{(2^i M)^{2+\delta}} < c_{11} \frac{\binom{n}{2}^2}{M^\delta}.$$

From (7) it follows that if  $M \geq c_{12}(\delta)$  then

$$\sum_{k: m_k \geq M} \binom{m_k}{2} < \frac{1}{2} \binom{n}{5} / (n-4).$$

Comparing it to (5) we conclude that

$$(8) \quad \left( \sum_{k: m_k < c_{12}(\delta)} 1 \right) \binom{c_{12}(\delta)}{2} > \frac{1}{2} \binom{n}{5} / (n-4).$$

Now (8) yields the desired lower bound

$$\text{cross}^* (\{P_1, \dots, P_n\}) = r \geq \sum_{k: m_k < c_{12}(\delta)} 1 > c_{13}(\delta) \binom{n}{4},$$

and Theorem 1.2 follows.  $\square$

ACKNOWLEDGEMENTS. The author wishes to thank Prof. P. Erdős and Prof. T. Zaslavsky (Ohio State University) for bringing the first problem to his attention. The author is very indebted to the referee for his corrections.

#### REFERENCES

- [1] BECK, J., On the "lattice property" of the plane and some problems of Dirac, Motzkin and Erdős in combinatorial geometry, *Combinatorica* 3 (1983), 281–297.
- [2] ERDŐS, P., On a problem of Grünbaum, *Canad. Math. Bull.* 15 (1972), 23–25. MR 47#5709.
- [3] GRÜNBAUM, B., *Arrangements and Spreads*, Regional Conference Series in Math. Nr. 10, Amer. Math. Soc., Providence, 1972. MR 46#6148.
- [4] KELLY, L. M. and MOSER, W. O. J., On the number of ordinary lines determined by  $n$  points, *Canad. J. Math.* 1 (1958), 210–219. MR 20#3494.
- [5] SZEMERÉDI, E. and TROTTER, W. T., Extremal problems in discrete geometry, *Combinatorica* 3 (1983) 381–392.
- [6] ZASLAVSKY, T., Extremal arrangements of hyperplanes (to appear in proceedings volume of the conference Discrete Geometry and Convexity Days, New York, April 2–3, 1982).

(Received June 20, 1983)

# ON THE ORDER OF CONVERGENCE OF A FINITE ELEMENT METHOD FOR THE BIHARMONIC EQUATION

L. VEIDINGER

Bramble and Zlámal considered in [1] a finite element, the so-called Argyris triangle for fourth order problems on a polygonal plane region  $R$  and gave asymptotic estimates for the order of convergence. They assumed, however, that the solution of the fourth order problem has square integrable fourth derivatives (at least) in  $R$ . This assumption is, in general, not satisfied even if  $R$  is convex (see [2], p. 288). In the present paper we shall give error bounds without this assumption.

1. Let  $R$  be a bounded open plane region whose boundary  $C$  consists of a finite number of simple closed polygons. Denote by  $A_\mu$  ( $\mu=1, 2, \dots, \nu$ ) the vertices of  $C$ .

We consider the Dirichlet problem for the biharmonic equation

$$\Delta^2 u(x, y) = f(x, y), \quad (x, y) \in R,$$

(1)

$$u|_C = \frac{\partial u}{\partial n}|_C = 0,$$

where  $n$  is the outward normal to  $C$ . We assume that the right-hand side  $f(x, y)$  is analytic in an open region  $G$  containing the closure of  $R$  in its interior.

Let  $k$  be a non-negative integer. We denote by  $W_2^{(k)}(R)$  the Hilbert space of all functions which together with their generalized derivatives up to the  $k$ -th order belong to  $L_2(R)$ . The norm is given by

$$\|v\|_{k,R}^2 = \sum_{j=0}^k |v|_{j,R}^2, \quad \text{where} \quad |v|_{j,R}^2 = \sum_{|i|=j} \|D^i v\|_{L_2(R)}^2.$$

Here and in the sequel we use the notation  $i=(i_1, i_2)$ ,  $|i|=i_1+i_2$ ,  $D^i v = \frac{\partial^{|i|} v}{\partial x^{i_1} \partial y^{i_2}}$ .

If  $s$  is a positive real number, then we define the Hilbert space  $W_2^{(s)}(R)$  with norm  $\|\cdot\|_{s,R}$  by Hilbert space interpolation between  $W_2^{(\lfloor s \rfloor)}(R)$  and  $W_2^{(\lfloor s \rfloor + 1)}(R)$  (see, for example, [3], p. 302).



It is well-known that under the above-mentioned assumptions the solution  $u(x, y)$  of the problem (1) minimizes the functional

$$(2) \quad F(v) = \frac{1}{2} \iint_R (\Delta v)^2 dx dy - \iint_R f v dx dy$$

in the space  $\dot{W}_2^{(2)}(R)$ . Here  $\dot{W}_2^{(2)}(R)$  is the space of functions which we get by completing in the norm  $\|\cdot\|_{2,R}$  the set of twice continuously differentiable functions with compact support in  $R$ .

LEMMA 1.  $u(x, y)$  is an analytic function of  $x$  as well as of  $y$  in  $R$ .

For a proof, see [4].

LEMMA 2.  $u(x, y)$  is analytic on  $C$ , excluding the vertices.

For a proof, see [5].

LEMMA 3. Let  $A_\mu = (x_{A_\mu}, y_{A_\mu})$  be a vertex of  $C$  with interior angle  $\alpha_\mu$  ( $0 < \alpha_\mu < 2\pi$ ). Let  $r_{A_\mu} = \sqrt{(x - x_{A_\mu})^2 + (y - y_{A_\mu})^2}$ . Let  $\bar{\alpha}_\mu = \operatorname{Re} z_0 + 1$  where  $z_0$  is the root with smallest positive real part of the (complex) equation

$$\sin^2 \alpha_\mu z - z^2 \sin^2 \alpha_\mu = 0.$$

Then

$$u(x, y) = u_1(x, y) + O(r_{A_\mu}^{\bar{\alpha}_\mu}),$$

where  $u_1(x, y)$  and its partial derivatives of all orders remain bounded when  $(x, y) \rightarrow A_\mu$  in  $R$ . This relation may be indefinitely formally differentiated.

This lemma follows from the results of Kondrat'ev (see [2]).

2. We triangulate  $R$ , i.e. we subdivide  $R$  into triangles such that any two triangles are either disjoint or have a common vertex or a common side. To every triangulation we associate two parameters:  $h, \vartheta$ .  $h$  is the largest side and  $\vartheta$  is the smallest angle of all triangles of the given triangulation. In the sequel we assume that

$$\vartheta \geq \vartheta_0 > 0,$$

where  $\vartheta_0$  does not depend on  $h$ . Denote by  $M_h$  the set of all triangles of the given triangulation.

Let  $P_1, P_2, P_3$  be the vertices of the triangle  $T \in M_h$  and let  $Q_1, Q_2, Q_3$  be the mid-points of the sides. We consider a polynomial of degree 5

$$p(x, y) = a_0 + a_1 x + \dots + a_{20} x y^4 + a_{21} y^5.$$

To determine such a polynomial we need 21 conditions. We choose them in the following way: we prescribe the values  $D^i p(P_j)$  ( $|i| \leq 2, j=1, 2, 3$ ) and the values  $\frac{\partial p(Q_j)}{\partial n}$  ( $j=1, 2, 3$ ), where  $n$  is the outward normal to the boundary of  $T$ . It can be proved (see [1], p. 810) that the polynomial  $p(x, y)$  is uniquely determined by these values.



We consider the values  $D^i p(P_j)$  and  $\frac{\partial p(Q_j)}{\partial n}$  ( $|i| \leq 2$ ,  $j=1, 2, 3$ ) at each interior node as parameters. Denote by  $\dot{H}_5(R)$  the class of functions defined on  $R$  which are equal on each triangle  $T \in M_h$  to the just introduced polynomial  $p(x, y)$  and satisfy the boundary conditions (see [6], p. 404). It can be shown (see [6], p. 404) that  $\dot{H}_5(R)$  is a finite dimensional subspace of  $\dot{W}_2^{(2)}(R)$ . We approximate the solution  $u(x, y)$  of the problem (1) by the function  $u_h(x, y)$  which minimizes the functional (2) in the space  $\dot{H}_5(R)$ .

3. THEOREM. Let  $u(x, y)$  be the solution of the problem (1) and let  $u_h(x, y)$  be the finite element approximation. Let  $\beta = \min_{\mu=1, 2, \dots, \nu} \bar{\alpha}_\mu$ . If  $\beta \leq 5$ , then

$$u(x, y) \in \dot{W}_2^{(1+\beta-\varepsilon)}(R) \cap \dot{W}_2^{(2)}(R)$$

and

$$(3) \quad \|u - u_h\|_{2,R} \leq c_1 h^{\beta-1-\varepsilon} \|u\|_{1+\beta-\varepsilon, R},$$

where  $\varepsilon$  is any positive real number and  $c_1$  is a positive constant which depends only on the right-hand side  $f(x, y)$  and the region  $R$ . If  $\beta > 5$ , then  $u(x, y) \in \dot{W}_2^{(6)}(R) \cap \dot{W}_2^{(2)}(R)$  and

$$(4) \quad \|u - u_h\|_{2,R} \leq c_2 h^4 \|u\|_{6,R},$$

where  $c_2$  is a positive constant which depends only on  $f(x, y)$  and  $R$ .

PROOF<sup>1</sup>. By Céa's lemma (see, for example, [7], p. 104) we have

$$(5) \quad \|u - u_h\|_{2,R} \leq c_3 \inf_{v \in \dot{H}_5(R)} \|u - v\|_{2,R}$$

where  $c_3$  is a positive constant which does not depend on  $h$ .

Denote by  $P$  the projection of  $\dot{W}_2^{(2)}(R)$  onto  $\dot{H}_5(R)$ . If  $z \in \dot{W}_2^{(2)}(R)$ , then

$$\inf_{v \in \dot{H}_5(R)} \|z - v\|_{2,R} = \|z - Pz\|_{2,R}$$

$P$  and  $T = I - P$  are continuous linear operators on  $\dot{W}_2^{(2)}(R)$ . In fact

$$(6) \quad \|Tz\|_{2,R} \leq \|z\|_{2,R} + \|Pz\|_{2,R} \leq 2\|z\|_{2,R}.$$

On the other hand, if  $z \in \dot{W}_2^{(2)}(R) \cap \dot{W}_2^{(6)}(R)$  then (see [1])

$$(7) \quad \|Tz\|_{2,R} \leq c_4 h^4 \|z\|_{6,R} \leq c_4 h^4 \|z\|_{6,R}.$$

(6) and (7) imply that  $T$  is a mapping from  $\dot{W}_2^{(2)}(R)$  into  $\dot{W}_2^{(2)}(R)$  with norm  $\leq 2$  and it is a mapping from  $\dot{W}_2^{(6)}(R) \cap \dot{W}_2^{(2)}(R)$  into  $\dot{W}_2^{(2)}(R)$  with norm  $\leq c_4 h^4$ . We have (see [8], p. 263) for  $\mu=1, 2, \dots, \nu$ ,  $r_{A_\mu}^\beta \in \dot{W}_2^{(\beta+1-\varepsilon)}(R)$ , where  $\varepsilon$  is any positive real number. Hence by the lemmas it follows that  $u(x, y) \in \dot{W}_2^{(\beta+1-\varepsilon)}(R) \cap \dot{W}_2^{(2)}(R)$ . If  $\beta > 5$ , then  $u(x, y) \in \dot{W}_2^{(6)}(R) \cap \dot{W}_2^{(2)}(R)$  and thus by (7) the inequality (4) imme-

<sup>1</sup> The fundamental idea of this proof goes back to Scott (see [3]).

diately follows. If  $\beta \leq 5$ , then by the so-called operator interpolation property (see, for example, [3], p. 301)  $T$  is a mapping from  $W_2^{(1+\beta-\varepsilon)}(R) \cap \tilde{W}_2^{(2)}(R)$  into  $\tilde{W}_2^{(2)}(R)$  with norm  $\leq c_5 h^{\beta-1-\varepsilon}$ , that is

$$\inf_{v \in \tilde{H}_0^2(R)} \|u - v\|_{2,R} = \|Tu\|_{2,R} \leq c_5 h^{\beta-1-\varepsilon} \|u\|_{1+\beta-\varepsilon, R}.$$

Substituting (8) into (5) we immediately obtain the inequality (3). This completes the proof of our Theorem.

Finally, we remark that our results can be transferred without any difficulty to all the finite elements considered in [7] on p. 355 and also to the curved elements considered in [9] (in the case when  $C$  is curved). Moreover, using the results of [10], the finite element methods for the biharmonic eigenvalue problem

$$\Delta^2 u(x, y) + \lambda u(x, y) = 0, \quad (x, y) \in R,$$

$$u \Big|_C = \frac{\partial u}{\partial n} \Big|_C = 0,$$

can be treated in the same way as those for the boundary value problem (1).

#### REFERENCES

- [1] BRAMBLE, J. H. and ZLÁMAL, M., Triangular elements in the finite element method, *Math. Comp.* **24** (1970), 809—820. *MR* **43**#8250.
- [2] KONDRAT'EV, V. A., Boundary value problems for elliptic equations with conical or angular points, *Trudy Moskov. Mat. Obšč.* **16** (1967), 209—292 (in Russian). *MR* **37**#1777.
- [3] SCOTT, R., *Applications of Banach Space interpolation to finite element theory*, *Functional analysis methods in numerical analysis*, ed. by M. Z. Nashed, Lect. Notes in Math., 701, Springer-Verlag, Berlin, 1979. *MR* **80b**: 65139.
- [4] JOHN, F., The fundamental solutions of linear elliptic differential equations with analytic coefficients, *Comm. Pure Appl. Math.* **3** (1950), 273—304. *MR* **31**—40.
- [5] MORREY, C. D. and NIRENBERG, L., On the analyticity of linear elliptic systems of partial differential equations, *Comm. Pure Appl. Math.* **10** (1957), 271—290. *MR* **19**—654.
- [6] ZLÁMAL, M., On the finite element method, *Numer. Math.* **12** (1968), 394—409. *MR* **39**#5074.
- [7] CIARLET, R. G., *The finite element method for elliptic problems*, North-Holland, Amsterdam, 1978. *MR* **58**#25001.
- [8] STRANG, G. and FIX, G., *An analysis of the finite element method*, Prentice-Hall, Inc., Englewood Cliffs, N. J., 1973. *MR* **56**#1747.
- [9] MANSFIELD, L., Approximation of the boundary in the finite element solution of fourth order problems, *SIAM J. Numer. Anal.* **15** (1978), 568—579.
- [10] VEIDINGER, L., On the order of convergence of finite element approximations to eigenvalues and eigenfunctions, *Period. Math. Hungar.* **8** (1977), 45—52.

(Received June 23, 1983)

MTA MATEMATIKAI KUTATÓ INTÉZETE  
P.O. BOX 127  
H—1364 BUDAPEST  
HUNGARY

TWO REMARKS ON THE PERMEABILITY OF LAYERS OF  
CONVEX BODIESA. FLORIAN and H. GROEMER<sup>1</sup>

We let  $E^d$  denote the  $d$ -dimensional Euclidean space. As usual, a *convex body* in  $E^d$  is defined as a compact convex subset of  $E^d$  with interior points. A collection  $\mathcal{D}$  of convex bodies whose interiors are mutually disjoint is called a *layer* in  $E^d$  if there are two parallel hyperplanes  $H_1, H_2$  in  $E^d$  such that all bodies from  $\mathcal{D}$  are between  $H_1$  and  $H_2$ . These hyperplanes are said to be *boundary planes* (if  $d=2$ , *boundary lines*) of  $\mathcal{D}$  if there are no two hyperplanes parallel to  $H_1, H_2$  which also include  $\mathcal{D}$  and have smaller mutual distance than  $H_1, H_2$ . The distance between the two boundary planes of  $\mathcal{D}$  is called the *width* of  $\mathcal{D}$ . We shall only deal with layers that satisfy the condition that every bounded subset of  $E^d$  intersects only a finite number of convex bodies of the layer. We concern ourselves with continuous rectifiable curves that lie entirely between the boundary planes  $H_1, H_2$  of the layer  $\mathcal{D}$ . Such a curve will be called a *path*, and we shall say that a path *crosses*  $\mathcal{D}$  if it has one endpoint in  $H_1$ , the other in  $H_2$ , and does not meet the interior of any convex body from  $\mathcal{D}$ . The length of a path  $\lambda$  will be denoted by  $l(\lambda)$ . Since we are interested in curves of small length that cross  $\mathcal{D}$  it will be convenient to define

$$L(\mathcal{D}) = \inf l(\lambda),$$

where the infimum is to be taken over all paths  $\lambda$  that cross  $\mathcal{D}$ . When we talk about circles, spheres,  $m$ -gons etc. we mean always the corresponding convex body, not just its boundary.

Following L. Fejes Tóth [2] we define the *permeability* of a layer  $\mathcal{D}$  by

$$p(\mathcal{D}) = \frac{w}{L(\mathcal{D})},$$

where  $w$  denotes the width of  $\mathcal{D}$ . For every layer  $\mathcal{D}$  we have obviously  $0 < p(\mathcal{D}) \leq 1$ . The permeabilities of layers of circles and some other plane convex bodies have been investigated in the articles [1] through [7]. It is shown in [2] that every layer of congruent circles has permeability greater than  $\sqrt{27}/2\pi = 0.8269\dots$ , and that there are layers of incongruent circles of permeability less than  $\sqrt{27}/2\pi$  (in fact, less than 0.8234). Hence, if  $\Gamma$  denotes the class of all layers of congruent circles and  $\Gamma^*$  the class of all layers of not necessarily congruent circles we have

$$\inf \{p(\mathcal{C}): \mathcal{C} \in \Gamma\} > \inf \{p(\mathcal{C}^*): \mathcal{C}^* \in \Gamma^*\}.$$

<sup>1</sup> Supported by National Science Foundation Research Grant MCS 8001578.

1980 *Mathematics Subject Classification*. Primary 52A45; Secondary 52A40.

*Key words and phrases*. Permeability, layers of convex bodies, regular polygons.

(We use frequently an asterisk to indicate that incongruent convex bodies may be involved.) In [3] L. Fejes Tóth has shown that parallelograms behave quite differently in this respect. If  $Q$  is a parallelogram and  $\Pi, \Pi^*$  denote, respectively, the classes of all layers of translates of  $Q$  and parallelograms similar to  $Q$ , then

$$\inf \{p(\mathcal{Q}): \mathcal{Q} \in \Pi\} = \inf \{p(\mathcal{Q}^*): \mathcal{Q}^* \in \Pi^*\}.$$

In Theorem 1 of the present paper we note that for  $m \geq 39$  the regular  $m$ -gon is of the same type as the circle. This leaves us with the interesting question to which type the regular  $m$ -gon belongs if  $m$  is one of the integers 3, 5, 6, ..., 38. For example, we do not know whether the regular hexagon is of the same type as the circle or the parallelogram.

Our second theorem provides a lower bound for the permeabilities of layers of convex bodies in  $E^d$ . In the case of congruent spheres this lower bound tends rapidly to 1 as  $d$  tends to infinity.

**THEOREM 1.** *For every  $m \geq 39$  there exists a layer  $\mathcal{P}^*$  of homothetic regular  $m$ -gons such that*

$$p(\mathcal{P}^*) < \inf p(\mathcal{P}),$$

where the infimum is to be taken over all layers  $\mathcal{P}$  of congruent regular  $m$ -gons.

**PROOF.** Let  $\mathcal{C}^*$  be a layer of width 1 of not necessarily congruent circles, and let  $\mathcal{P}^*$  be a layer of regular  $m$ -gons that are inscribed in the circles of  $\mathcal{C}^*$  and mutually homothetic. (One of the  $m$ -gons may be inscribed arbitrarily, all the others are then uniquely determined. Until the very end of the proof we assume only  $m \geq 3$ .)

For any  $\varepsilon > 0$  there exists a path  $\beta_1$  that crosses  $\mathcal{P}^*$  and satisfies the inequality

$$(1) \quad l(\beta_1) < L(\mathcal{P}^*) + \varepsilon.$$

Let  $H_1, H_2$  denote the boundary lines of  $\mathcal{C}^*$ . Since the boundary lines of  $\mathcal{P}^*$  are between  $H_1$  and  $H_2$  we obviously may supplement the path  $\beta_1$  by two line segments orthogonal to  $H_1, H_2$  to obtain a path  $\beta_2$  that connects  $H_1$  and  $H_2$ . If  $\mathcal{P}^*$  has width  $w$  it follows that  $w \leq 1$  and

$$(2) \quad l(\beta_2) = l(\beta_1) + 1 - w.$$

We now modify the path  $\beta_2$  so that we obtain a path, say  $\beta_3$ , that crosses  $\mathcal{C}^*$  and is, in a certain sense, not much longer than  $\beta_2$ . If  $\beta_2$  does not meet the interior of any circle from  $\mathcal{C}^*$  we set  $\beta_3 = \beta_2$ . Let us now assume that there is a  $C \in \mathcal{C}^*$  such that  $\beta_2$  meets  $\text{int } C$ . Then, part of  $\beta_2$  is contained in one of the  $m$  open convex regions, say  $S$ , of  $(\text{int } C) \setminus P$ , where  $P$  is the regular  $m$ -gon from  $\mathcal{P}^*$  inscribed in  $C$  (see Fig. 1). If  $\beta_2(t)$  is the parametric representation of  $\beta_2$  there is a point  $\beta_2(a)$  on the circular boundary of  $S$  where  $\beta_2$  enters  $S$ , and a point  $\beta_2(b)$  where it exits. Let  $\lambda$  be the subarc of  $\beta_2$  defined by  $a \leq t \leq b$ , and  $\lambda'$  the subarc on the circular boundary between  $\beta_2(a)$  and  $\beta_2(b)$ . If  $R$  is the radius of  $C$  and  $\alpha$  the angle  $l(\lambda')/R$  we have  $\alpha \leq 2\pi/m$ ,  $l(\lambda) \geq 2R \sin \frac{\alpha}{2}$ , and therefore

$$(3) \quad \frac{l(\lambda')}{l(\lambda)} \leq \frac{\pi/m}{\sin(\pi/m)}.$$

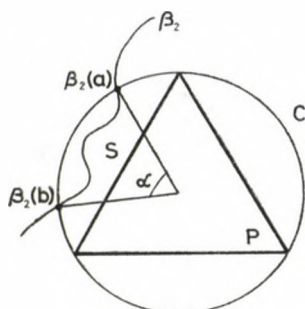


Fig. 1

We now remove from  $\beta_2$  the subarc  $\lambda$  and replace it by  $\lambda'$ , and apply a corresponding replacement to any other subarc of  $\beta_2$  that meets a region of the same kind as  $S$ . The result is a path  $\beta_3$  that crosses  $\mathcal{C}^*$  and, because of (3), has the property that

$$l(\beta_3) \leq \frac{\pi/m}{\sin(\pi/m)} l(\beta_2).$$

If this inequality is combined with (1) and (2) we obtain

$$l(\beta_3) < \frac{\pi/m}{\sin(\pi/m)} \left( 1 - w + \frac{w}{p(\mathcal{P}^*)} + \varepsilon \right).$$

Because of  $w \leq 1$  and  $p(\mathcal{P}^*) \leq 1$  we have  $1 - w + w/p(\mathcal{P}^*) \leq 1/p(\mathcal{P}^*)$ , and it follows that

$$l(\beta_3) < \frac{\pi/m}{\sin(\pi/m)} \left( \frac{1}{p(\mathcal{P}^*)} + \varepsilon \right)$$

and therefore

$$L(\mathcal{C}^*) < \frac{\pi/m}{\sin(\pi/m)} \left( \frac{1}{p(\mathcal{P}^*)} + \varepsilon \right).$$

Since the width of  $\mathcal{C}^*$  was assumed to be 1 we have  $p(\mathcal{C}^*) = 1/L(\mathcal{C}^*)$ , and since  $\varepsilon > 0$  was arbitrary we finally obtain

$$(4) \quad p(\mathcal{P}^*) \leq \frac{\pi/m}{\sin(\pi/m)} p(\mathcal{C}^*).$$

Let now  $\mathcal{P}$  be a layer of width 1 of congruent regular  $m$ -gons, and let  $\mathcal{C}$  be the layer consisting of the circles inscribed in the  $m$ -gons of  $\mathcal{P}$ . If  $\varepsilon > 0$  is given, there is a path  $\gamma_1$  that crosses  $\mathcal{C}$  and satisfies the inequality

$$(5) \quad l(\gamma_1) < L(\mathcal{C}) + \varepsilon.$$

Similarly as before, we supplement  $\gamma_1$  by two line segments orthogonal to the boundary lines of  $\mathcal{P}$  so that the resulting path, say  $\gamma_2$ , connects these two lines. If  $w$  is the width of  $\mathcal{C}$  we have  $w \leq 1$  and

$$(6) \quad l(\gamma_2) = l(\gamma_1) + 1 - w.$$

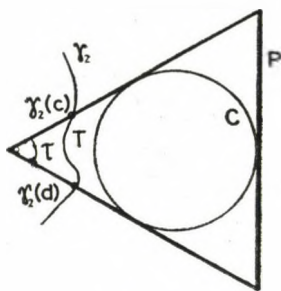


Fig. 2

Our aim is now to modify  $\gamma_2$  so that the new path, say  $\gamma_3$ , crosses  $\mathcal{P}$  and is not much longer than  $\gamma_2$ . If  $\gamma_2$  does not meet any of the interiors of the  $m$ -gons from  $\mathcal{P}$  we simply set  $\gamma_3 = \gamma_2$ . If this is not the case,  $\gamma_2$  meets the interior of some  $P \in \mathcal{P}$  and passes therefore through one of the open regions, let us call it  $T$ , bounded by two adjacent sides of  $P$  and a circular arc from the corresponding inscribed circle (see Fig. 2). Let  $\gamma_2(c)$  be the point on the polygonal boundary of  $T$  where  $\gamma_2$  enters  $T$ , and  $\gamma_2(d)$  the point where  $\gamma_2$  exits  $T$ . If  $\mu$  is the subarc of  $\gamma_2$  defined by the parametric restriction  $c \leq t \leq d$ , and  $\mu'$  denotes the polygonal boundary of  $T$  between  $\gamma_2(c)$  and  $\gamma_2(d)$  we find

$$(7) \quad \frac{l(\mu')}{l(\mu)} \leq \frac{1}{\cos(\pi/m)}.$$

Indeed, if  $\mu'$  consists of only one line segment this is completely obvious, and if it consists of two line segments one obtains (7) as an immediate consequence of the fact that for any triangle with side lengths  $x, y, z$  and respective angles  $\varrho, \sigma, \tau = \pi - \frac{2\pi}{m}$  we have

$$\frac{x+y}{z} = \frac{\cos((\varrho-\sigma)/2)}{\sin(\tau/2)} \leq \frac{1}{\cos(\pi/m)}.$$

We now replace the subarc  $\mu$  of  $\gamma_2$  by  $\mu'$  and carry out the corresponding replacements for all other subarcs of  $\gamma_2$  that meet a region of the same kind as  $T$ . This yields a curve  $\gamma_3$  that crosses  $\mathcal{P}$  and, by (7), satisfies the inequality

$$l(\gamma_3) \leq \frac{1}{\cos(\pi/m)} l(\gamma_2).$$

If this inequality is combined with (5) and (6) we obtain

$$l(\gamma_3) < \frac{1}{\cos(\pi/m)} (1 - w + L(\mathcal{C}) + \varepsilon).$$

Similarly as in the discussion of  $\mathcal{C}^*$  and  $\mathcal{P}^*$ , we may deduce that

$$L(\mathcal{P}) \leq \frac{1}{\cos(\pi/m)} \frac{1}{p(\mathcal{C})}$$



and therefore

$$(8) \quad p(\mathcal{C}) \equiv \frac{1}{\cos(\pi/m)} p(\mathcal{P}).$$

From (4) and (8) it follows that for any layer  $\mathcal{C}^*$  of possibly incongruent circles, and any layer  $\mathcal{P}$  of congruent regular  $m$ -gons there exist a layer  $\mathcal{P}^*$  of homothetic regular  $m$ -gons and a layer  $\mathcal{C}$  of congruent circles so that

$$(9) \quad \frac{p(\mathcal{P}^*)}{p(\mathcal{P})} \equiv \frac{2\pi/m}{\sin(2\pi/m)} \frac{p(\mathcal{C}^*)}{p(\mathcal{C})}.$$

(The assumption that  $\mathcal{C}^*$  and  $\mathcal{P}$  have width 1 is obviously immaterial.)

We now apply (9) to a layer  $\mathcal{C}^*$  of incongruent circles with  $p(\mathcal{C}^*) < 0.8234$ . As already remarked, the existence of such a layer has been shown by L. Fejes Tóth [2] who has also proved that for any layer  $\mathcal{C}$  of congruent circles  $p(\mathcal{C}) > \sqrt{27}/2\pi$ . Thus, no matter how  $\mathcal{P}$  is chosen,

$$\frac{p(\mathcal{C}^*)}{p(\mathcal{C})} < 0.995655.$$

Using this inequality and the easily established fact that for  $m \geq 39$

$$\frac{2\pi/m}{\sin(2\pi/m)} < 1.00434,$$

it follows from (9) that

$$p(\mathcal{P}^*) < 0.999977 p(\mathcal{P}).$$

Since  $\mathcal{P}$  can be any layer of congruent regular  $m$ -gons the proof of Theorem 1 is finished.

To formulate our second theorem we need the following definition. If  $\mathcal{K} = \{K_1, K_2, \dots\}$  is a layer in  $E^d$  we may introduce a cartesian coordinate system in  $E^d$  so that the two boundary planes of  $\mathcal{K}$  consist of all points  $(x_1, \dots, x_d)$  with  $x_d = 0$  and  $x_d = w$ , respectively. For any  $r > 0$  we let  $Z(r)$  denote the cylinder defined by  $x_1^2 + \dots + x_{d-1}^2 \leq r^2$ ,  $0 \leq x_d \leq w$ . Then, we define the upper density of  $\mathcal{K}$  by

$$(10) \quad \varrho(\mathcal{K}) = \lim_{r \rightarrow \infty} \frac{1}{v(Z(r))} \sum_{K_i \cap Z(r) \neq \emptyset} v(K_i),$$

where  $v$  denotes the volume in  $E^d$ . It is not difficult to show that  $\varrho(\mathcal{K})$  does not depend on the particular coordinate system introduced for the purpose of this definition. We say that a line is orthogonal to  $\mathcal{K}$  if it is orthogonal to the boundary planes of  $\mathcal{K}$ .  $|\cdot|$  denotes the Euclidean norm in  $E^d$ .

**THEOREM 2.** Let  $\mathcal{K}$  be a layer in  $E^d$  of permeability  $p$  and upper density  $\varrho$ . Furthermore, assume that there exists a constant  $\gamma$  with the following property: For any  $K \in \mathcal{K}$  and any line that is orthogonal to  $\mathcal{K}$  and intersects the boundary of  $K$  in two points, say  $s, t$ , there is a path on the boundary of  $K$  that connects  $s$  and  $t$ , and has length at most  $\gamma|s-t|$ . Then,

$$(11) \quad p \geq \frac{1}{1 + (\gamma - 1)\varrho}.$$



PROOF. We use the coordinate system and the cylinder  $Z(r)$  that have been described in connection with the definition (10). For any point  $b=(x_1, \dots, x_{n-1}, 0)$  in the base, say  $B(r)$ , of  $Z(r)$  we consider the line  $G(x_1, \dots, x_{n-1})$  that passes through  $b$  and is orthogonal to  $\mathcal{K}=\{K_1, K_2, \dots\}$ . Let  $f(x_1, \dots, x_{n-1})$  be the linear measure of  $G(x_1, \dots, x_{n-1}) \cap (\bigcup_i K_i)$ . Then, we find

$$(12) \quad \int_{B(r)} f(x_1, \dots, x_{n-1}) dx_1 \dots dx_{n-1} = v(Z(r) \cap (\bigcup_i K_i)) \equiv \sum_{K_i \cap Z(r) \neq \emptyset} v(K_i).$$

If for some  $i$  the interior of the body  $K_i$  meets the line  $G(x_1, \dots, x_{n-1})$  we replace the line segment  $G(x_1, \dots, x_{n-1}) \cap K_i$  by a path on bdr  $K_i$  that connects the end-points  $s, t$  of this line segment and has length at most  $\gamma|s-t|$ . If we perform corresponding replacements for all other bodies  $K_j$  with  $(\text{int } K_j) \cap G(x_1, \dots, x_{n-1}) \neq \emptyset$  we obtain a path that crosses  $\mathcal{K}$  and has length at most

$$w - f(x_1, \dots, x_{n-1}) + \gamma f(x_1, \dots, x_{n-1}),$$

where  $w$  denotes again the width of  $\mathcal{K}$ . It follows immediately that

$$L(\mathcal{K}) \leq w + (\gamma - 1)f(x_1, \dots, x_{n-1}).$$

If both sides of this inequality are integrated over  $B(r)$ , and if we apply (12) and observe that  $B(r)$  has  $((d-1)$ -dimensional) volume  $v(Z(r))/w$ , we obtain

$$L(\mathcal{K}) \frac{v(Z(r))}{w} \leq v(Z(r)) + (\gamma - 1) \sum_{K_i \cap Z(r) \neq \emptyset} v(K_i).$$

(11) is now an obvious consequence of this inequality and (10).

To illustrate the applicability of Theorem 2 let us consider some special cases. Assume first that  $\mathcal{K}$  is a layer of (not necessarily congruent) equilateral triangles with one side parallel to the boundary lines of  $\mathcal{K}$ . We have  $\varrho \leq 1$  and may take  $\gamma = \sqrt{3}$ . Then Theorem 2 yields  $p \geq 1/\sqrt{3}$ , and this is actually best possible since a layer as shown in Fig. 3 is easily seen to have permeability arbitrarily close to  $1/\sqrt{3}$  if the width of the layer is sufficiently large. If equilateral triangles in any position are permitted one can take  $\gamma = 2$  to obtain the estimate  $p \geq 1/2$ .

As another example let  $\mathcal{K}$  be a layer of congruent spheres in  $E^d$ , and let  $p_d$  denote the infimum of all permeabilities of layers of congruent spheres. Since  $\varrho(\mathcal{K})$  cannot be larger than the upper density of packings of congruent spheres in  $E^d$ , and since one may obviously take  $\gamma = \frac{\pi}{2}$ , we obtain from Theorem 2 that

$$(13) \quad \frac{1}{p_d} \leq 1 + \left(\frac{\pi}{2} - 1\right) \sigma_d,$$

where  $\sigma_d$  is the upper bound of Rogers [9] for the density of sphere packings in  $E^d$ . Thus, for  $d \rightarrow \infty$  we have  $\sigma_d \sim \frac{d}{e} 2^{-d/2}$  and therefore

$$\lim_{d \rightarrow \infty} p_d = 1.$$

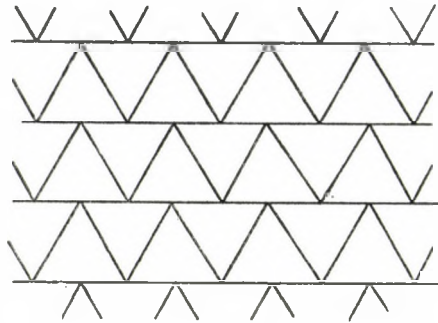


Fig. 3

If one uses instead of  $\sigma_d$  the density estimates of Kabatyanskii and Levenshtein [8] one obtains the inequality

$$\frac{1}{p_d} - 1 \leq 2^{-(0.599 + o(1))d} \quad (d \rightarrow \infty)$$

which is stronger than (13) if  $d$  is large enough.

We finally remark that in the proof of Theorem 2 the assumption of convexity of the sets in  $\mathcal{H}$  was practical but not essential. Hence, substantial generalizations would be possible.

## REFERENCES

- [1] BOLLOBÁS, B., Remarks to a paper of L. Fejes Tóth, *Studia Sci. Math. Hungar.* **3** (1968), 373—379. *MR* 39#866.
- [2] FEJES TÓTH, L., On the permeability of a circle-layer, *Studia Sci. Math. Hungar.* **1** (1966), 5—10. *MR* 34#6640.
- [3] FEJES TÓTH, L., On the permeability of a layer of parallelograms, *Studia Sci. Math. Hungar.* **3** (1968), 195—200. *MR* 38#611.
- [4] FLORIAN, A., On the permeability of layers of discs, *Studia Sci. Math. Hungar.* **13** (1978), 125—132. *MR* 83a: 52014.
- [5] FLORIAN, A., Über die Durchlässigkeit gewisser Scheibenschichten, *Österreich. Akad. Wiss. Math.—Naturw. Kl. Sitzungsber. II* **188** (1979), 417—427. *MR* 82c: 52004.
- [6] FLORIAN, A., Über die Durchlässigkeit einer Schicht konvexer Scheiben, *Studia Sci. Math. Hungar.* **15** (1980), 201—213. *MR* 84h: 52020.
- [7] HORTOBÁGYI, I., Über die Durchlässigkeit einer aus Scheiben konstanter Breite bestehenden Schicht, *Studia Sci. Math. Hungar.* **11** (1976), 383—387. *MR* 81c: 52014.
- [8] KABATYANSKII, G. A. and LEVENSSTEIN, V. I., Bounds for packings on a sphere and in space, *Problemy Peredachi Informacii* **14** (1978), 3—25 (Russian). English translation in: *Problems of Information Transmission* **14** (1978), 1—17. *MR* 58#24018.
- [9] ROGERS, C. A., The packing of equal spheres, *Proc. London Math. Soc.* (3) **8** (1958), 609—620. *MR* 21#847.

(Received June 3, 1983)

INSTITUT FÜR MATHEMATIK  
UNIVERSITÄT SALZBURG  
HELLBRUNNERSTRASSE 34  
A-5020 SALZBURG  
AUSTRIA

DEPARTMENT OF MATHEMATICS  
THE UNIVERSITY OF ARIZONA  
TUCSON, AZ 85721  
U.S.A.



# SUR LES SURFACES À GROUPES CONTINUS $G_2$ DE SIMILITUDE PROJECTIVES EN ELLES-MÊMES ET SUR LES SURFACES COMPLEXES

FROIM MARCUS

*Hommage à la mémoire de O. Mayer*

1. La notion de correspondance projective-conforme et similaire entre deux surfaces non-quadriques que l'on doit à Bompiani [2] a été généralisée à une surface et sa transformée infinitésimale en elle-même, dans les notes [3] et [4]. Une transformation infinitésimale asymptotique  $(t \cdot i)$  d'une surface  $x(u, v)$  en elle-même est représentée par le symbole  $X = k \frac{\partial}{\partial u} + l \frac{\partial}{\partial v}$ ,  $k = k(u)$ ,  $l = l(v)$  ou par les équations

$$(1.1) \quad \bar{u} = u + \varepsilon k, \quad \bar{v} = v + \varepsilon l,$$

$\varepsilon$  un paramètre indépendant de  $u, v$  dont on néglige le carré.  $X$  sera une  $t \cdot i$  projective-conforme  $(t \cdot i \cdot p-c)$  si l'on peut trouver une fonction

$$(1.2) \quad \varrho = 1 + \varepsilon \sigma \quad (\sigma = \sigma(u, v) \neq \text{const.})$$

telle que

$$(1.3) \quad \frac{\bar{\beta} d\bar{u}^3 + \bar{\gamma} d\bar{v}^3}{2 d\bar{u} d\bar{v}} = \varrho \frac{\beta du^3 + \gamma dv^3}{2 du dv},$$

soit vérifiée au premier ordre où  $\frac{\beta du^3 + \gamma dv^3}{2 du dv}$  est l'élément linéaire de Fubini [1] de  $(x)$  et  $\varrho$  le module dilatation de l'élément.

$X$  est une similitude  $(t \cdot i \cdot p-s)$  si  $\sigma = \text{const.} \neq 0$  et une déformation projective de Fubini  $(d \cdot i \cdot p)$  si  $\sigma = 0$  (les collinéations exclues). Les résultats de [3] et [4] ont été étendus par beaucoup de résultats par O. Mayer et F. Marcus dans [5], dont nous rappelons en particulier les suivants:

(1) Si une surface non réglée admet un  $G_2$  continu de  $t \cdot p-c$ , elle admet un  $G_3$  et elle est isothermo asymptotique ou surface  $(F)$  de Fubini.

(2) Une surface  $(F)$  possède un groupe continue maximum  $G_3$  de  $t \cdot p-c$  en elle-même.

2. Le paragraphe 7 de [5] traite le problème de déterminer les surfaces qui possèdent un groupe  $G_2$  de similitudes projectives  $(t \cdot p-s)$ . On démontre qu'en dehors des surfaces de coïncidence qui possèdent un  $G_3$  il existent deux et seulement deux éléments linéaires dont les surfaces correspondantes possèdent un  $G_2$  maximum

de  $t \cdot p$ -s en elles-mêmes engendré respectivement par les  $t \cdot i$

$$X_1 = \frac{\partial}{\partial u}, X_2 = \frac{\partial}{\partial v} \quad \text{et par} \quad X_1 = b \frac{\partial}{\partial u} - a \frac{\partial}{\partial v} \quad (ab \neq 0), \quad X_2 = u \frac{\partial}{\partial u} + v \frac{\partial}{\partial v}.$$

Les deux types d'éléments linéaires sont:

$$(2.1) \quad \beta = \gamma = Ce^{\omega} \quad (\omega = au + bv \neq 0),$$

et

$$(2.2) \quad \beta = \gamma = C\omega^{\tau-1} \quad (\tau \neq 1),$$

qui sont réalisables et à qui correspondent six classes de surfaces, la classe (I) pour l'élément linéaires (2.1) et les autres pour l'élément (2.2). Parmi les six classes, les surfaces de la classe (II) sont les plus intéressantes. Des conditions d'intégrabilité [1] on trouve si  $ab \neq 0$

$$(2.3) \quad L = -\frac{3}{2} C^2 \frac{a}{b} \omega^{2\tau-2} + U_1(u), \quad M = -\frac{3}{2} C^2 \frac{b}{a} \omega^{2\tau-2} + V_1(v),$$

et si  $\tau=0$  et  $C^2=ab$  on est conduit à la classe de surfaces déterminées (formellement) par

$$\beta = \gamma = \sqrt{ab}(au+bv)^{-1} \quad (ab \neq 0).$$

$$(II) \quad L = -\frac{3}{2} a^2 (au+bv)^{-2} + C' au^2 + C'' u + C''' b$$

$$M = -\frac{3b^2}{2} (au+bv)^2 + C' bv^2 - C'' v + C''' a.$$

On observe dans [5] p. 405 qu'elles dépendent seulement de deux constantes essentielles car on peut y faire

$$1: C' = \pm 1, C'' = 0; \quad 2: C' = 0, C'' = 1, C''' = 0;$$

(2.4)

$$3: C' = C'' = 0, C''' \pm 1; \quad 4: C' = C'' = C''' = 0.$$

Au quatre groupes des constantes (2.4) correspondent évidemment quatre sous-classes de surfaces projectivement applicables entre elles, mais avec propriétés différentes, spécialement en relation avec les congruences —  $W$  ayant une quadrique comme nappe focale et qui n'ont pas été considérées dans [5]. L'objet de cette note est de compléter nos résultats de [5].

3. De (2.2) l'on tire

$$(3.1) \quad \frac{\partial^2 \log \beta}{\partial u \partial v} = \beta^2, \quad (\gamma = \beta)$$

c'est-à-dire que les asymptotiques des deux familles sur les surfaces correspondantes appartiennent à des complex linéaires. Donc elles sont des surfaces de Terracini [6]

ou encore d'après G. Bol [7] surfaces complexes et qui sont de trois espèces; celles d'espèce (III) ont été découvert par Wilczynski [8].

D'après Fubini—Čech [1] paragraphe C, p. 115, on a généralement

$$\beta = \frac{\sqrt{U'V'}}{U+V} = \gamma,$$

$$(3.2) \quad L = -\frac{\partial^2 \log \beta}{\partial u^2} - \frac{1}{2} \left( \frac{\partial \log \beta}{\partial u} \right)^2 + U_1,$$

$$M = -\frac{\partial^2 \log \beta}{\partial v^2} - \frac{1}{2} \left( \frac{\partial \log \beta}{\partial v} \right)^2 + V_1,$$

où

$$(3.3) \quad U_1 = \frac{kU^2 + (l-h)U + p}{U'}, \quad V_1 = \frac{kV^2 + (l-h)V + p}{V'},$$

$k, l, h, p$  sont des constantes pas toutes essentielles. En observant que dans notre cas  $U=au, V=bv$ , on aura

$$(3.4) \quad C' = k, \quad C'' = l-h, \quad C''' = \frac{p}{ab},$$

et de (2.4) l'on tire

$$(3.5) \quad \begin{aligned} 1^\circ \quad k = \pm 1, \quad l = h; & \quad 2^\circ \quad k = 0, \quad l-h = 1; \\ (\mathcal{P}) \quad 3^\circ \quad k = 0, \quad l = h, \quad p = \pm ab; & \quad 4^\circ \quad k = l = h = p = 0. \end{aligned}$$

On aura donc pour les quatre sous-classes:

$$(3.6) \quad \beta = \sqrt{ab} (au+bv)^{-1} = \gamma$$

$$\begin{aligned} L_1 &= \alpha \pm au^2 + \frac{p}{a}; & M_1 &= \delta \pm bv^2 + \frac{p}{b}; \\ (\mathcal{P}') \quad L_2 &= \alpha + u; & M_2 &= \delta - v; \\ L_3 &= \alpha \pm ab; & M_3 &= \delta \pm ab; \\ L_4 &= \alpha; & M_4 &= \delta \end{aligned}$$

où

$$(3.6) \quad \alpha = -\frac{3}{2} a^2 (au+bv)^{-2}, \quad \delta = -\frac{3}{2} b^2 (au+bv)^{-2}$$

et les première trois classes de surfaces sont projectivement applicables sur la dernière qui est une surface de Wilczynski.

4. Dans le paragraphe 47 C p. 272—275 Fubini démontre le Théorème: Si  $S$  est une surface dont chacune des asymptotiques appartiennent à un complex linéaire, alors  $S$  est une nappe focale de congruences  $W$  dont la seconde nappe  $S$  est une quadrique. (On observe aussi qu'il y a des exceptions.)

Soit  $S$  une telle surface. Alors pour que  $\bar{S}$  soit une quadrique on doit avoir (lieu cité)

$$(4.1) \quad N = P \frac{(HV+K)(K-HU)}{U+V}, \quad (P = \text{const.} \neq 0)$$

$U$  et  $V$  sont solutions de (3.1) et

$$(4.2) \quad H(h-l) - 2kK = -\frac{P}{2}H, \quad K(h-l) - 2pH = \frac{P}{2}K,$$

$H, K$  constantes non simultanément nulles et

$$(4.3) \quad kK^2 - pH^2 = -\frac{P}{2}KH.$$

Par élimination de  $H$  et  $K$  l'on tire

$$(4.4) \quad (h-l)^2 - 4kp = \frac{P^2}{2};$$

que déterminera  $P$  (sauf pour le signe), et pour chaque valeur de  $P$  on obtient  $H$  et  $K$ . Observons que  $N$  est la bien connue fonction de la méthode générale de Fubini [1] pour l'étude des congruences— $W$ . Si  $N = \text{const.} \neq 0$ , la congruence appartient à un complexe linéaire, si  $N = 0$ , le complexe est spécial et la seconde nappe focale dégénère. C'est précisément le cas d'exception. Donc si  $N \neq \text{const.}$  alors chacune des surfaces  $S$  est une nappe focale de deux congruences  $W$ , ayant une quadrique pour seconde nappe focale. Observons encore que les deux congruences sont engendrés par les tangentes aux courbes

$$(4.5) \quad \frac{du}{A} = \frac{dv}{B'}$$

avec

$$(4.5') \quad A = \frac{\sqrt{U'}}{U+V}(HV+K), \quad B' = \frac{\sqrt{V'}}{U+V}(-HU+K).$$

De (4.4) et pour les constantes du groupe 1° de ( $\mathcal{P}$ ) résultera

$$(4.6) \quad P^2 = \mp 16p \neq 0, \quad PK + 4pH = 0.$$

Observant que  $U=au, V=bv$ , on aura (4.7)  $N=N(u, v) \neq \text{const.}$  Donc les surfaces respectives sont nappes focales de deux congruences— $W$  dont les secondes nappes sont des quadriques régulières. Elles sont donc surfaces complexes d'espèce (I) de Terracini.

Les congruences sont engendrés par les tangentes aux courbes

$$(4.8) \quad \frac{du}{\sqrt{a}(Hbv+K)} = \frac{dv}{\sqrt{b}(K-Hau)}.$$

Pour le deuxième groupe de ( $\mathcal{P}$ ) on aura

$$(4.9) \quad P = \pm 2.$$

Pour  $P=2$ , il résulte  $H \neq 0, K=0$ , et si  $P=-2, H=0, K \neq 0$ . On aura respec-



tivement

$$(4.10) \quad N_1 = -\frac{2H^2 ab uv}{au + bv}, \quad N_2 = -\frac{2K^2}{au + bv}.$$

Les surfaces sont aussi d'espèce (I) et l'on a aussi

$$(4.11) \quad (\log N)_{uv} = (\log \beta)_{uv}.$$

Les congruences  $W$  sont engendrés par les tangentes aux courbes

$$(4.12) \quad bdu^2 - adv^2 = 0.$$

Passons maintenant au troisième groupe. De (4.4) et  $(\mathcal{P})$  résultera

$$(4.13) \quad P = 0,$$

et par conséquence

$$(4.14) \quad N = 0.$$

Observant que  $p = \pm ab \neq 0$ , l'on tire de (4.2)

$$(4.15) \quad H = 0, \quad K \neq 0, \quad A = \frac{\sqrt{a} K}{au + bv}, \quad B = \frac{\sqrt{b} K}{au + bv}.$$

*La congruence— $W$  appartient à un complexe linéaire spéciale — dont la deuxième nappe focale est dégénéré et se réduit à l'axe de complexe. Donc les surfaces qui correspondent au troisième groupe de  $(\mathcal{P})$  font exception au théorème énoncé ci-dessus. Enfin considérons le quatrième groupe des constantes  $(\mathcal{P})$ . De (4.4) l'on tire*

$$(4.16) \quad P = 0, \quad N = 0,$$

*donc les surfaces correspondant font aussi exception au théorème. Elles sont des surfaces limite de Tzitzeica—Wilczynski [1], c'est-à-dire la première directrice de Wilczynski passe par un point fixe et la deuxième directrice est contenu dans un plan fixe passant par le point commun des premières directrices. Mais elles sont aussi des surfaces de Demoulin car les sommets du quadrilatère de Demoulin sont confondues. Elles sont aussi minima-projective, car d'après O. Mayer [9] il résulte de la dernière de  $(\mathcal{P}')$  compte tenue de (3.5) et (3.6)*

$$(4.17) \quad \beta M_v + 2M\beta_v + \beta_{vvv} = 0, \quad (\beta L_u + 2L\beta_u + \beta_{uuu} = 0).$$

On observe de même que la condition (4.17) n'est pas remplie par les autres invariants  $(L_1, M_1)(L_2, M_2)$  et  $(L_3, M_3)$  d'accord avec la proposition suivante démontré dans [10]: *Parmi les surfaces complexes seules les surfaces de troisième espèce sont minima-projectives.* Enfin observant que (4.4) ne dépend pas de  $U$  et  $V$ , on peut appliquer les résultats aux surfaces complexes les plus générales en utilisant le groupe  $(\mathcal{P})$ .

5. Nous voulons maintenant donner la représentation paramétrique de la surface complexe qui correspond au groupe 4° de  $(\mathcal{P})$ , c'est-à-dire pour laquelle

$$(5.1) \quad \beta = \sqrt{ab}(au+bv)^{-1},$$

$$L = -\frac{3}{2}a^2(au+bv)^{-2}, \quad M = -\frac{3}{2}b^2(au+bv)^{-2}.$$

Nous ne pouvons pas l'obtenir de

$$(5.2) \quad x_1 = (u-v)(U'-V')-2(U+V), \quad x_2 = U'-V', \quad x_3 = u+v, \quad x_4 = 1,$$

$U, V$  sont des fonctions arbitraires de  $y$  et  $v$  respectivement non réductible à des polynômes quadratiques, que donne la représentation paramétrique des surfaces d'espèce (III). Voir [11].

Mais on peut utiliser le résultat suivant de Wilczynski [8] p. 157—160:

*En coordonnées asymptotiques, les équations*

$$(5.3) \quad \frac{y_1}{2\sqrt{\psi}} = U_2V_3 + U_3V_2 + \int \frac{U_2 - UU_3}{\sqrt{U'}} du + \int \frac{V_2 - VV_3}{\sqrt{V'}} dv; \quad \frac{y_2}{\sqrt{\psi}} = U_2 - V_2;$$

$$\frac{y_3}{\sqrt{\psi}} = U_3 + V_3; \quad \frac{y_4}{\sqrt{\psi}} = 1; \quad (U(u), V(v) \text{ sont fonctions arbitraires})$$

où

$$(5.4) \quad \psi = \frac{\sqrt{U'V'}}{2(U+V)}; \quad U_2 = \int \frac{U du}{\sqrt{U'}}, \quad U_3 = \int \frac{du}{\sqrt{U'}},$$

$$V_2 = \int \frac{V dv}{\sqrt{V'}}, \quad V_3 = \int \frac{dv}{\sqrt{V'}}$$

déterminent les  $y_i$  comme coordonnées homogènes des surfaces complexes non réglées dont les courbes directrices sont indéterminées:

En observant que  $U=au$ ,  $V=bv$ , et  $2\psi=\beta$ , l'on tire

$$(5.5) \quad U_2 = \sqrt{a} \frac{u^2}{2}, \quad U_3 = \frac{u}{\sqrt{a}}, \quad V_2 = \sqrt{b} \frac{v^2}{2}, \quad V_3 = \frac{v}{\sqrt{b}}.$$

Par conséquent

$$(5.6) \quad y_1 = \frac{\sqrt{a} u^2 v + \sqrt{b} u v^2}{\sqrt{ab}} - \frac{u^3 + v^3}{3}; \quad y_2 = \frac{\sqrt{a} u^2 - \sqrt{b} v^2}{\sqrt{ab}};$$

$$y_3 = \frac{\sqrt{b} u + \sqrt{a} v}{\sqrt{ab}}; \quad y_4 = 1,$$

est la représentation paramétrique de la surface (5.1). Les coordonnées (5.6) sont les

solutions du système

$$(5.7) \quad \xi_{uu} = \theta_u \xi_u - \beta \xi_v, \quad \xi_{vv} = -\beta \xi_u + \theta_v \xi_v,$$

où

$$\theta_u = -(\log \beta)_u, \quad \theta_v = -(\log \beta)_v, \quad \beta = \frac{\sqrt{ab}}{au + bv}.$$

Elles sont donc les coordonnées tangentielles de la surface.

#### RÉFÉRENCES

- [1] FUBINI, G. et ČECH, E., *Geometria proiettiva differenziale*, Bologna, 1926.
- [2] BOMPIANI, E., I fondamenti geometrici della teoria proiettiva delle superficie. Appendice in [1].
- [3] MARCUS, F., Asupra suprafețelor ce admit deformatii infinitesimale proiectiv-simile, *Acad. R. P. Romîne Fil. Iași Stud. Cerc. Sti. Mat.* **12** (1961), 291—314. *MR* **26** # 6879.
- [4] MARCUS, F., Asupra deformărilor infinitesimale proiectiv-conforme etc., *Ibid.* **13** (1962), 109—128. *MR* **26** # 6880.
- [5] MAYER, O. et MARCUS, F., Surfaces à groupes continus de transformations projectives-conformes en elles-mêmes, *An. Ști. Univ. "Al. I. Cuza" Iași Sect. Mat.* **9** (1963), 387—408. *MR* **32** # 4619.
- [6] TERRACINI, A., Sulle superficie aventi un sistema o entrambi di asintotiche in complessi lineari, Appendice (IV) in [1], 771—782.
- [7] BOL, G., *Projektive Differentialgeometrie*, 2. Teil, Göttingen, 1954, 318—319. *MR* **16**—1150.
- [8] WILCZYŃSKI, E. J., Über Flächen mit unbestimmten Direktrixkurven, *Math. Ann.* **76** (1915), 122—160.
- [9] MAYER, O., Contribution à l'étude des surfaces minima-projectives, *Bull. des Sciences Mathématiques*, 1932.
- [10] MARCUS, F., Sur les surfaces de troisième espèce de Terracini, *Czechoslovak Math. J.* **6**(81) (1956), 559—562. *MR* **19**—879.
- [11] MARCUS, F., Some remarks on non-ruled surfaces whose asymptotes, etc., *An. Ști. Univ. "Al. I. Cuza" Iași* **27** (1981).

(Received July 1, 1983)

DEPARTMENT OF MATHEMATICS  
TECHNION  
ISRAEL INSTITUTE OF TECHNOLOGY  
HAIFA  
ISRAEL



# ON THE EIGENFUNCTIONS OF FIRST- AND SECOND-ORDER DIFFERENTIAL OPERATORS

V. KOMORNIK

Let  $G \subset \mathbb{R}$  be a bounded open interval,  $n \in \mathbb{N}$ ,  $p \in [1, \infty]$ ,  $q_1, \dots, q_n \in L^p(G)$  arbitrary complex functions and consider the formal differential operator

$$(1) \quad Lu = u^{(n)} + q_1 u^{(n-1)} + \dots + q_n u.$$

Given a complex number  $\lambda$ , the function  $u: \bar{G} \rightarrow \mathbb{C}$ ,  $u \equiv 0$  is called an eigenfunction of order  $-1$  of the operator  $L$  with the eigenvalue  $\lambda$ . More generally, as in usual, a function  $u: \bar{G} \rightarrow \mathbb{C}$ ,  $u \not\equiv 0$  is called an eigenfunction of order  $m$  ( $m=0, 1, \dots$ ) of the operator  $L$  with the eigenvalue  $\lambda$  if the following two conditions are satisfied:

—  $u, u', \dots, u^{(n-1)}$  are absolute continuous on  $\bar{G}$ ;

— there exists an eigenfunction  $u^*$  of order  $m-1$  of the operator  $L$  with the eigenvalue  $\lambda$  such that

$$(2) \quad (Lu)(x) = \lambda u(x) + u^*(x) \quad \text{a.e. on } \bar{G}.$$

Developing the results of the papers [2] and [4], we have recently proved the following estimates (see [7]):

“Given any eigenfunction  $u$  of order  $m$  of the operator  $L$  with some eigenvalue  $\lambda = \mu^n$ , the following estimates are valid:

$$(3) \quad \|u^{(i)}\|_{L^{p'}(G)} \leq C_m (1 + |\mu|)^i \|u\|_{L^{p'}(G)} \quad (i = 0, \dots, n-1),$$

$$(4) \quad \|u^*\|_{L^{p'}(G)} \leq C_m (1 + |\lambda|) \|u\|_{L^{p'}(G)},$$

$$(5) \quad \|u\|_{L^\infty(G)} \leq C_m (1 + |\mu|)^{1/q} \|u\|_{L^q(G)} \quad (q \in [1, \infty]).$$

Furthermore, for  $|\lambda|$  sufficiently large

$$(6) \quad \|u^{(i)}\|_{L^{p'}(G)} \leq B_m (1 + |\mu|)^i \|u\|_{L^{p'}(G)} \quad (i = 0, \dots, n-1),$$

$$(7) \quad \|u^{(i)}\|_{L^\infty(G)} \leq C_m (1 + |\mu|)^{1/q} \|u^{(i)}\|_{L^q(G)} \quad (i = 1, 2, \dots, n-1, q \in [1, \infty])$$

( $C_m, B_m$  are positive constants).”

One can see easily that these estimates cannot be improved if  $n \geq 3$ . However, in case  $n \leq 2$  the estimates (4), (5), (7) are not optimal. (Concerning (4) and (5) the exact estimates are proved in [2].) The aim of this paper is to find the optimal estimates in these cases, too. We shall improve and generalize the results of the papers [1], [2], [3], [5].

1980 *Mathematics Subject Classification*. Primary 34C11; Secondary 15A18.

*Key words and phrases*. Eigenfunction of higher order.

Throughout this paper let  $u$  be an eigenfunction of order  $m$  of the operator  $L$  with the eigenvalue  $\lambda = \mu^n$  and let us introduce the continuous functions  $u_j$ ,  $j \leq m$  by the formulas

$$(8) \quad u_m = u \quad \text{and} \quad u_{j-1} = Lu_j - \lambda u_j \quad \text{a.e. on } \bar{G}.$$

### 1. First-order operators

In this section we consider the case  $n=1$ . We shall prove the following three theorems:

**THEOREM 1.** *There exists a constant  $C_m^*$  such that*

$$(9) \quad \|u^*\|_{L^{p'}(G)} \leq C_m^* (1 + |\operatorname{Re} \lambda|) \|u\|_{L^p(G)},$$

$$(10) \quad \|u\|_{L^\infty(G)} \leq C_m^* (1 + |\operatorname{Re} \lambda|)^{1/q} \|u\|_{L^q(G)} \quad (q \in [1, \infty]).$$

Let us now introduce the notations  $G=(a, b)$  and

$$d_1(x) = \begin{cases} x-a & \text{if } \operatorname{Re} \lambda < 0, \\ b-x & \text{if } \operatorname{Re} \lambda \geq 0. \end{cases}$$

**THEOREM 2.** *In case  $\|q_1\|_{L^1(G)} < 1$  there exists a constant  $D_m^*$  such that*

$$(11) \quad \sup_{x \in G} |u(x)(1 + |\operatorname{Re} \lambda| d_1(x))^{-m} \exp(|\operatorname{Re} \lambda| d_1(x))| \leq D_m^* \|u\|_{L^\infty(G)}.$$

*In the general case there exist constants  $\alpha \in (0, 1)$  and  $E_m^*$  such that*

$$(12) \quad \sup_{x \in G} |u(x) \exp(\alpha |\operatorname{Re} \lambda| d_1(x))| \leq E_m^* \|u\|_{L^\infty(G)}.$$

**THEOREM 3.** *There exists a positive constant  $B_m^*$  such that*

$$(13) \quad \|u\|_{L^\infty(G)} \leq B_m^* (1 + |\operatorname{Re} \lambda|)^{1/q} \|u\|_{L^q(G)} \quad (q \in [1, \infty]).$$

In what follows, we shall systematically use the following formula, being a special case of [4], Theorem 1 (however, this very simple form is due to I. Joó, see [8]):

$$(14) \quad \begin{aligned} & t^j u_{m-j}(x) = \\ &= \sum_{k=1}^{m+1} c_{mjk} e^{-k\lambda t} \left[ u_m(x+kt) + \sum_{r=0}^m \int_x^{x+kt} \frac{(x+kt-\tau)^r}{r!} e^{\lambda(x+kt-\tau)} q_1(\tau) u_{m-r}(\tau) d\tau \right] \end{aligned}$$

whenever  $x \in \bar{G}$ ,  $x+(m+1)t \in \bar{G}$ ,  $j=0, 1, \dots$ . Here the numbers  $c_{mjk}$  are absolute constants. For brevity, we set  $N:=m+1$ . Hence follows the existence of a constant  $A_m$  such that

$$(15) \quad |t^j u_{m-j}(x)| \leq A_m \sum_{k=1}^N |u_m(x+kt)| + A_m \sum_{r=0}^m |t|^r \|q_1\|_{L^p(x, x+Nt)} \|u_{m-r}\|_{L^{p'}(x, x+Nt)}$$

whenever  $x \in \bar{G}$ ,  $x+Nt \in \bar{G}$ ,  $|\operatorname{Re} \lambda t| \leq 1$ ,  $j \in \{0, 1, \dots, m\}$ .

PROOF OF THEOREM 1. Let us set

$$|G| = b - a,$$

$$R = \min \left\{ \frac{|G|}{2N}, |\operatorname{Re} \lambda|^{-1} \right\} \quad (0^{-1} := \infty),$$

$$M_{p'} = \max \{R^j \|u_{m-j}\|_{L^{p'}(G)} : j = 0, 1, \dots, m\}.$$

Applying (15) with

$$x \in \left[ a, \frac{a+b}{2} \right], \quad t = R, \quad j \in \{0, 1, \dots, m\},$$

$$R^j |u_{m-j}(x)| \leq A_m \sum_{k=1}^N |u_m(x+kR)| + NA_m \|q_1\|_{L^p(G)} M_{p'}$$

whence

$$R^j \|u_{m-j}\|_{L^{p'}(a, (a+b)/2)} \leq NA_m \|u_m\|_{L^{p'}(G)} + NA_m \|q_1\|_{L^p(G)} |G|^{1/p'} M_{p'}.$$

The same estimate can be obtained for  $\|u_{m-j}\|_{L^{p'}(\frac{a+b}{2}, b)}$ , too, hence

$$M_{p'} \leq 2NA_m \|u_m\|_{L^{p'}(G)} + 2NA_m |G|^{1/p'} \|q_1\|_{L^p(G)} M_{p'}.$$

If

$$(16) \quad 4NA_m |G|^{1/p'} \|q_1\|_{L^p(G)} \leq 1$$

then we can conclude

$$M_{p'} \leq 4NA_m \|u_m\|_{L^{p'}(G)}$$

whence (9) follows.

In the general case we can divide  $G$  to the union of finitely many subintervals  $G_i$  such that for each index  $i$

$$4NA_m |G_i|^{1/p'} \|q_1\|_{L^p(G_i)} \leq 1;$$

then

$$\|u^*\|_{L^{p'}(G_i)} \leq C_{m,i}^* (1 + |\operatorname{Re} \lambda|) \|u\|_{L^{p'}(G_i)}$$

for each  $i$ , and (9) follows with  $C_m^* = \sum_i C_{m,i}^*$ .

To prove (10), let us apply (15) with  $x \in \left[ a, \frac{a+b}{2} \right]$ ,  $t \in (0, R)$  and  $j \in \{0, 1, \dots, m\}$  (we can put  $p=1$ ):

$$t^j |u_{m-j}(x)| \leq A_m \sum_{k=1}^N |u_m(x+kt)| + NA_m \|q_1\|_{L^1(G)} M_\infty.$$

Applying the transformation  $NR^{-1} \int_0^R dt$  and using the Hölder inequality,

$$R^j \|u_{m-j}\|_{L^\infty(a, (a+b)/2)} \leq N^2 A_m R^{-1/q} \|u_m\|_{L^q(G)} + N^2 A_m \|q_1\|_{L^1(G)} M_\infty$$

whence by symmetry

$$M_\infty \leq N^2 A_m R^{-1/q} \|u_m\|_{L^q(G)} + N^2 A_m \|q_1\|_{L^1(G)} M_\infty.$$



If

$$(17) \quad 2N^2 A_m \|q_1\|_{L^1(G)} \leq 1$$

then

$$M_\infty \leq 2N^2 A_m R^{-1/2} \|u_m\|_{L^2(G)}$$

which implies (10). The general case hence follows as before.

The theorem is proved.  $\square$

**PROOF OF THEOREM 2.** It suffices to consider the case  $Q := \|q_1\|_{L^1(G)} < 1$ , the general case hence follows by the method seen in the preceding proof. We work by induction on  $m$ . For  $m = -1$  the assertion is obvious. Let now  $u$  be an eigenfunction of order  $m$  and suppose the theorem is true for eigenfunctions of order  $< m$ . One can easily see that

$$(18) \quad u(x)e^{\lambda t} = u(x+t) + \int_x^{x+t} e^{\lambda(x+t-\tau)} [q_1(\tau)u(\tau) - u^*(\tau)] d\tau$$

whenever  $x \in \bar{G}$  and  $x+t \in \bar{G}$ . Suppose that  $\operatorname{Re} \lambda \geq 0$ ; the case  $\operatorname{Re} \lambda < 0$  is analogous. Applying (18) for any  $x \in G$  with  $t = d_1(x) = b - x$ ,

$$(19) \quad \begin{aligned} & |u(x)e^{\operatorname{Re} \lambda(b-x)}| \leq \\ & \leq |u(b)| + Q \sup_{\tau \in (x,b)} |u(\tau)e^{\operatorname{Re} \lambda(b-\tau)}| + (b-x) \sup_{\tau \in (x,b)} |u^*(\tau)e^{\operatorname{Re} \lambda(b-\tau)}|. \end{aligned}$$

Now using the inductive hypothesis and the estimate (9),

$$(20) \quad \begin{aligned} & (b-x) \sup_{\tau \in (x,b)} |u^*(\tau)e^{\operatorname{Re} \lambda(b-\tau)}| \leq \\ & \leq (b-x)(1 + \operatorname{Re} \lambda(b-x))^{m-1} \sup_{\tau \in (x,b)} |u^*(\tau)e^{\operatorname{Re} \lambda(b-\tau)}(1 + \operatorname{Re} \lambda(b-\tau))^{1-m}| \leq \\ & \leq (b-x)(1 + \operatorname{Re} \lambda(b-x))^{m-1} D_{m-1}^* \|u^*\|_{L^\infty(G)} \leq \\ & \leq (b-x)(1 + \operatorname{Re} \lambda(b-x))^{m-1} D_{m-1}^* C_m^* (1 + \operatorname{Re} \lambda) \|u\|_{L^\infty(G)} \leq \\ & \leq (1 + \operatorname{Re} \lambda(b-x))^m (1 + |G|) D_{m-1}^* C_m^* \|u\|_{L^\infty(G)}. \end{aligned}$$

(19) and (20) yield

$$\begin{aligned} & \sup_{x \in \bar{G}} |u(x)(1 + \operatorname{Re} \lambda(b-x))^{-m} e^{\operatorname{Re} \lambda(b-x)}| \leq \\ & \leq Q \sup_{x \in \bar{G}} |u(x)(1 + \operatorname{Re} \lambda(b-x))^{-m} e^{\operatorname{Re} \lambda(b-x)}| + (1 + (1 + |G|) D_{m-1}^* C_m^*) \|u\|_{L^\infty(G)}; \end{aligned}$$

hence (11) follows because  $Q < 1$  and the proof is finished.  $\square$

**PROOF OF THEOREM 3.** It is an easy consequence of the estimate (12) (see also [3]).  $\square$

## 2. Second-order operators

In this section we turn to the case  $n=2$ . The following three theorems will be proved:

THEOREM 4. *There exists a constant  $C_m^*$  such that*

$$(21) \quad \|u^*\|_{L^{p'}(G)} \leq C_m^*(1+|\mu|)(1+|\operatorname{Re} \mu|) \|u\|_{L^{p'}(G)},$$

$$(22) \quad \|u\|_{L^\infty(G)} \leq C_m^*(1+|\operatorname{Re} \mu|)^{1/q} \|u\|_{L^q(G)} \quad (q \in [1, \infty]);$$

furthermore for  $|\lambda|$  sufficiently large

$$(23) \quad \|u'\|_{L^\infty(G)} \leq C_m^*(1+|\operatorname{Re} \mu|)^{1/q} \|u'\|_{L^q(G)} \quad (q \in [1, \infty]).$$

Let us now introduce the notations  $G=(a, b)$  and

$$d_2(x) = \operatorname{dist}(x, \partial G) \quad (= \min \{x-a, b-x\}) \quad (x \in G).$$

THEOREM 5. *In case  $\|q_s\|_{L^1(G)} < 2^{-m-1}$ ,  $s=1, 2$ , there exists a constant  $D_m^*$  such that*

$$(24) \quad \sup_{x \in G} |u^{(i)}(x)(1+|\operatorname{Re} \mu|d_2(x))^{-m} \exp(|\operatorname{Re} \mu|d_2(x))| \leq D_m^* \|u^{(i)}\|_{L^\infty(G)} \quad (i=0, 1).$$

In the general case there exist constants  $\alpha \in (0, 1)$  and  $E_m^*$  such that

$$(25) \quad \sup_{x \in G} |u^{(i)}(x) \exp(\alpha |\operatorname{Re} \mu| d_2(x))| \leq E_m^* \|u^{(i)}\|_{L^\infty(G)} \quad (i=0, 1).$$

THEOREM 6. *There exists a positive constant  $B_m^*$  such that*

$$(26) \quad \|u^{(i)}\|_{L^\infty(G)} \leq B_m^*(1+|\operatorname{Re} \mu|)^{1/q} \|u^{(i)}\|_{L^q(G)} \quad (q \in [1, \infty], i=0, 1).$$

REMARKS. The estimates (21), (22) were proved in [2] for the case  $q_1 \equiv 0$ . The following proof uses the ideas of the paper [5]. The estimate (23) seems to be new even for the case  $q_1 \equiv 0$ , too.

The Theorems 5 and 6 were proved in [3] for the case  $i=0$  and  $q_1 \equiv 0$ .

Let us set for  $r \in \{0, 1, \dots\}$   $K_{2,r}: \mathbb{C} \times \mathbb{R} \rightarrow \mathbb{C}$  with

$$K_{2,0}(\mu, x) = \frac{\operatorname{sh} \mu x}{\mu},$$

$$K_{2,r}(\mu, x) = \int_0^x K_{2,0}(\mu, x-t) K_{2,r-1}(\mu, t) dt \quad (r > 0).$$

We shall use the following formula and estimates (see [4], [5], [7]):

There exist holomorphic integral functions  $f_{2m}$  and  $f_{2mjik} = f_{2,m,j,i,k}$  ( $j \in \{0, \dots, m\}$ ,  $i \in \{0, 1\}$ ,  $i_0 \in \{0, i\}$ ,  $k \in \{1, \dots, 2m+2 =: N\}$ ) such that

$$(27) \quad f_{2m}(\mu t) t^{2j+i-i_0} u_{m-j}^{(i)}(x) = \sum_{k=1}^N f_{2,m,j,i-i_0,k}(\mu t) [u_m^{(i_0)}(x+kt) + \\ + \sum_{r=0}^m \sum_{s=1}^2 \int_x^{x+kt} D_2^{i_0} K_{2,r}(\mu, x+kt-\tau) q_s(\tau) u_{m-r}^{(2-s)}(\tau) d\tau]$$

whenever  $x \in \bar{G}$  and  $x+Nt \in \bar{G}$ . Furthermore, the following estimates hold true with some absolute constant  $A_m^*$ :

$$(28) \quad \text{For any } z \in \mathbb{C} \text{ with } |z| \geq 1 \text{ and } |\operatorname{Re} z| \leq 1 \text{ there exists } \alpha \in [1/2, 1] \text{ such that } f_{2m}(\alpha z) \neq 0 \text{ and}$$

$$|f_{2mjik}(\alpha z)| \leq A_m^* |z|^{j+i} |f_{2m}(\alpha z)| \quad (j \in \{0, \dots, m\}, i \in \{0, 1\}, k \in \{1, \dots, N\}).$$

For  $j=i-i_0=0$  the formula (27) can be simplified by  $f_{2m}(\mu t)$  and for any  $z \in \mathbb{C}$  with  $|\operatorname{Re} z| \leq 1$ ,

$$(29) \quad |f_{2m00k}(z)| \leq A_m^* |f_{2m}(z)| \quad (k \in \{1, \dots, N\}).$$

$$(30) \quad |D_2^{(i_0)} K_{2,r}(\mu, y)| \leq A_m^* \left( \frac{1+|\mu y|}{2} \right)^r |\mu|^{-2r-1+i_0}$$

$$\text{if} \quad |\operatorname{Re} \mu y| \leq N \quad (r \in \{0, \dots, m\}, i_0 \in \{0, 1\}).$$

REMARKS. Using the explicit (but a little complicated) expressions for the functions  $f_{2m}$ ,  $f_{2mjik}$ ,  $K_{2,r}$  in [4], one can easily see that

- $f_{2m}(z)$  is a constant ( $\neq 0$ ) multiple of  $\left(\frac{\operatorname{sh} z}{z}\right)^{(m+1)^2}$ ;
- $f_{2mjik}(z)$  is a linear combination of some functions of type

$$z^{-(m+1)^2} z^{2j+i-s} P_{2mjik}(e^z, e^{-z}) \quad (s \in \{j, \dots, \min\{2j+i, m\}\})$$

with some polynomials of two variables  $P_{2mjik}$ ;

- $K_{2,r}(\mu, y)$  has the form

$$\mu^{-2r-1} (P_{2,r}(\mu y) e^{\mu y} + Q_{2,r}(\mu y) e^{-\mu y})$$

where  $P_{2,r}$  and  $Q_{2,r}$  are polynomials of degree  $\leq r$ . These properties imply (28) and (30). (29) is proved in [4].

In the sequel we shall use the notations  $G = |b-a|$  and

$$Q = \max \{ \|q_1\|_{L^p(G)}, \|q_2\|_{L^p(G)} \}.$$

PROOF OF THEOREM 4. It suffices to consider the case when  $|G|$  and  $Q$  are small: the general case hence follows as in Section 1. Moreover, it suffices to deal with  $|\lambda|$  sufficiently large. Indeed, for  $|\lambda|$  bounded (21)–(22) follow from (4)–(5).

Setting

$$R = \min \left\{ \frac{|G|}{2N}, |\operatorname{Re} \mu|^{-1} \right\} \quad (0^{-1} := \infty),$$

$$M_{p'} = \max \{ R^j |\mu|^{-j-i} \|u_{m-j}^{(i)}\|_{L^{p'}(G)} : j \in \{0, \dots, m\}, i \in \{0, 1\} \}$$

and applying (27), (28) and (30) with  $x \in \left[ a, \frac{a+b}{2} \right]$ ,  $t=R$ ,  $j \in \{0, \dots, m\}$ ,  $i \in \{0, 1\}$ , there exists an absolute constant  $A_m$  and  $R_0 \in \left[ \frac{R}{2}, R \right]$  (independently of  $x$ ) such that

$$\begin{aligned} R^{2j+i} |\mu^{(i)}_{m-j}(x)| &\leq A_m |\mu R|^{j+i} \sum_{k=1}^N |u_m(x+kR_0)| + \\ &+ NA_m |\mu R|^{j+i} \sum_{r=0}^m \sum_{s=1}^2 |\mu R|^r |\mu|^{-2r-1} Q \|u_{m-r}^{(2-s)}\|_{L^{p'}(G)} \end{aligned}$$

whenever

$$|\mu| \geq \frac{2N}{|G|} \quad (\text{i.e. } |\mu R| \geq 1).$$

Let us now suppose that

$$(31) \quad |\mu| \geq \max \left\{ \frac{2N}{|G|}, 1 \right\};$$

then hence we obtain

$$R^j |\mu|^{-j-i} \|u_{m-j}^{(i)}\|_{L^{p'}(a, (a+b)/2)} \leq NA_m \|u_m\|_{L^{p'}(G)} + N^2 A_m |G|^{1/p'} Q M_{p'}.$$

and by symmetry

$$M_{p'} \leq 2NA_m \|u_m\|_{L^{p'}(G)} + 2N^2 A_m |G|^{1/p'} Q M_{p'}.$$

Under the hypothesis

$$(32) \quad 4N^2 A_m |G|^{1/p'} Q \leq 1$$

we can conclude

$$(33) \quad M_{p'} \leq 4NA_m \|u_m\|_{L^{p'}(G)}$$

which proves (21) (and (3) again for  $n=2$ ).

To prove (22) and (23), we can assume that  $p=1$  (and  $p'=\infty$ ). Applying (27), (29) and (30) with  $x \in \left[ a, \frac{a+b}{2} \right]$ ,  $t \in (0, R)$ ,  $j \in \{0, \dots, m\}$  and  $i \in \{0, 1\}$ , we obtain with an absolute constant  $A_m$  the inequality

$$\begin{aligned} |\mu^{-i} u_m^{(i)}(x)| &\leq \\ &\leq A_m \sum_{k=1}^N |\mu^{-i} u_m^{(i)}(x+kt)| + NA_m \sum_{r=0}^m \sum_{s=1}^2 |\mu R|^r |\mu|^{-2r-1} Q \|u_{m-r}^{(2-s)}\|_{L^\infty(G)} \end{aligned}$$

whence

$$|\mu^{-i} u_m^{(i)}(x)| \leq A_m \sum_{k=1}^N |\mu^{-i} u_m^{(i)}(x+kt)| + N^2 A_m Q M_\infty.$$

if  $|\mu| \geq 1$ . Using the transformation  $R^{-1} \int_0^R dt$ ,

$$|\mu|^{-i} \|u_m^{(i)}(x)\| \leq NA_m |\mu|^{-i} R^{-1/q} \|u_m^{(i)}\|_{L^q(G)} + N^2 A_m Q M_\infty,$$

and by symmetry

$$|\mu|^{-i} \|u_m^{(i)}\|_{L^\infty(G)} \leq NA_m |\mu|^{-i} R^{-1/q} \|u_m^{(i)}\|_{L^q(G)} + N^2 A_m Q M_\infty.$$

Assuming now (31) and (32) satisfied, we can apply (33); we obtain

$$(34) \quad |\mu|^{-i} \|u_m^{(i)}\|_{L^\infty(G)} \leq NA_m |\mu|^{-i} R^{-1/q} \|u_m^{(i)}\|_{L^q(G)} + 4N^3 A_m^2 Q \|u_m\|_{L^\infty(G)}.$$

In case  $i=0$  hence we obtain under the hypothesis

$$(35) \quad 8N^3 A_m^2 Q \leq 1$$

the estimates

$$(36) \quad \|u_m\|_{L^\infty(G)} \leq 2NA_m R^{-1/q} \|u_m\|_{L^q(G)}$$

and (22) is proved.

In case  $i=1$  for  $|\lambda|$  sufficiently large (6) and (34) yield

$$\|u'_m\|_{L^\infty(G)} \leq NA_m R^{-1/q} \|u'_m\|_{L^q(G)} + 4N^3 A_m^2 B_m^{-1} Q \|u'_m\|_{L^\infty(G)}$$

whence (23) follows if

$$(37) \quad 8N^3 A_m^2 B_m^{-1} Q \leq 1.$$

The theorem is proved.  $\square$

PROOF OF THEOREM 5. We consider only the case

$$(38) \quad Q < 2^{-m-1};$$

the general case hence follows easily.

The assertion is obvious if  $m = -1$ . Working by induction, assume that  $m \geq 0$  and that the theorem is true for all the eigenfunctions of order  $< m$ . Obviously, it suffices to show (24) for  $|\lambda|$  sufficiently large.

One can easily see (cf. [2]) that

$$(39) \quad \begin{aligned} 2u(x) \operatorname{ch}(\mu t) &= u(x-t) + u(x+t) + \\ &+ \int_{x-t}^{x+t} \frac{\operatorname{sh} \mu(t-|x-\xi|)}{\mu} [q_1(\xi) u'(\xi) + q_2(\xi) u(\xi) - u^*(\xi)] d\xi, \end{aligned}$$

$$(40) \quad \begin{aligned} 2\mu^{-1} u'(x) \operatorname{ch}(\mu t) &= \mu^{-1} u'(x-t) + \mu^{-1} u'(x+t) + \\ &+ \int_{x-t}^{x+t} \operatorname{sgn}(\xi-x) \frac{\operatorname{ch} \mu(t-|x-\xi|)}{\mu} [q_1(\xi) u'(\xi) + q_2(\xi) u(\xi) - u^*(\xi)] d\xi \end{aligned}$$

whenever  $x-t \in \bar{G}$  and  $x+t \in \bar{G}$ . Let us set for brevity

$$(41) \quad \begin{aligned} M &= \sup_{x \in G} |\mu(x)(1 + |\operatorname{Re} \mu| d_2(x))^{-m} \exp(|\operatorname{Re} \mu| d_2(x))|, \\ M' &= \sup_{x \in G} |\mu^{-1} u'(x)(1 + |\operatorname{Re} \mu| d_2(x))^{-m} \exp(|\operatorname{Re} \mu| d_2(x))|, \\ M^* &= \sup_{x \in G} |u^*(x)(1 + |\operatorname{Re} \mu| d_2(x))^{1-m} \exp(|\operatorname{Re} \mu| d_2(x))|. \end{aligned}$$

If  $|\mu| \equiv 1$  then putting  $t = d_2(x)$  in (39)–(40) and taking into account that  $d_2(\xi) \leq 2d_2(x)$  if  $|x - \xi| \leq d_2(x)$ ,

$$(42) \quad |u^{-i} u^{(i)}(x) \exp(|\operatorname{Re} \mu| d_2(x))| \leq 2|\mu|^{-i} \|u^{(i)}\|_{L^\infty(G)} + 2^m Q(1 + |\operatorname{Re} \mu| d_2(x))^m (M + M') + 2d_2(x) |\mu|^{-1} (1 + |\operatorname{Re} \mu| d_2(x))^{m-1} M^* \quad (i = 0, 1).$$

Applying (21) and the induction hypothesis, for  $|\mu| \equiv 1$  we obtain

$$(43) \quad \begin{aligned} d_2(x) |\mu|^{-1} M^* &\leq D_{m-1}^* d_2(x) |\mu|^{-1} \|u^*\|_{L^\infty(G)} \leq \\ &\leq C_m^* D_{m-1}^* \frac{1 + |\mu|}{|\mu|} d_2(x) (1 + |\operatorname{Re} \mu|) \|u\|_{L^\infty(G)} \leq \\ &\leq 2C_m^* D_{m-1}^* (1 + |G|) (1 + |\operatorname{Re} \mu| d_2(x)) \|u\|_{L^\infty(G)}. \end{aligned}$$

Now (42) and (43) imply

$$\begin{aligned} M &\leq 2\|u\|_{L^\infty(G)} + 4C_m^* D_{m-1}^* (1 + |G|) \|u\|_{L^\infty(G)} + 2^m Q(M + M'), \\ M' &\leq 2|\mu|^{-1} \|u'\|_{L^\infty(G)} + 4C_m^* D_{m-1}^* (1 + |G|) \|u\|_{L^\infty(G)} + 2^m Q(M + M'). \end{aligned}$$

Hence in view of (3)

$$M + M' \leq [2 + 4C_m + 8C_m^* D_{m-1}^* (1 + |G|)] \|u\|_{L^\infty(G)} + 2^{m+1} Q(M + M')$$

and taking into account (38),

$$(44) \quad M + M' \leq (1 - 2^{m+1} Q)^{-1} [2 + 4C_m + 8C_m^* D_{m-1}^* (1 + |G|)] \|u\|_{L^\infty(G)}.$$

The case  $i=0$  of (24) hence follows at once. The case  $i=1$  of (24) follows from (44) and (6).

The theorem is proved.  $\square$

PROOF OF THEOREM 6. This theorem is an immediate consequence of the estimate (25) in Theorem 5.  $\square$

#### REFERENCES

- [1] Ильин, В. А. и Йо, И., Равномерная оценка собственных функций и оценка сверху числа собственных значений оператора Штурма—Лиувилля с потенциалом из класса  $L^p$ , *Differencialnye Uravnenija* 15 (1979), 1164—1174. MR 80i: 34033.
- [2] Joó, I., Upper estimates for the eigenfunctions of the Schrödinger operator, *Acta Sci. Math. (Szeged)* 44 (1982), 87—93.
- [3] KOMORNIK, V., Lower estimates for the eigenfunctions of the Schrödinger operator, *Acta Sci. Math. (Szeged)* 44 (1982), 95—98.
- [4] KOMORNIK, V., Upper estimates for the eigenfunctions of higher order of a linear differential operator, *Acta Sci. Math. (Szeged)* 45 (1983), 261—271.

- [5] KOMORNIK, V., Generalization of a theorem of Joó, *Ann. Univ. Sci. Budapest. Eötvös Sect. Math.* **27** (1985), 59—64.
- [6] Joó, I. and KOMORNIK, V., On the equiconvergence of expansions by Riesz bases formed by eigenfunctions of the Schrödinger operator, *Acta Sci. Math. (Szeged)* **46** (1983), 357—375.
- [7] KOMORNIK, V., Some new estimates for the eigenfunctions of higher order of a linear differential operator, *Acta Math. Acad. Sci. Hungar.* **45** (1985), 451—457.
- [8] Joó, I., Remarks to a paper of V. Komornik, *Acta Sci. Math. (Szeged)* **47** (1984), 201—204.

( Received July 17, 1983 )

EÖTVÖS LORÁND TUDOMÁNYEGYETEM  
TERMÉSZETTUDOMÁNYI KAR  
ANALÍZIS TANSZÉK  
MÚZEUM KRT. 6—8.  
H—1088 BUDAPEST  
HUNGARY



# DIFFERENCES OF THE TERMS OF LINEAR RECURRENCES

PÉTER KISS

A linear recurrence  $U = \{U_n\}_{n=0}^\infty$  of order  $k$  is defined by the recursion

$$U_n = a_1 U_{n-1} + a_2 U_{n-2} + \dots + a_k U_{n-k},$$

where the initial terms  $U_0, U_1, \dots, U_{k-1}$  and the  $a_i$ 's ( $i=1, 2, \dots, k$ ) are given rational integers. We can suppose that the initial terms are not all zero and  $a_k \neq 0$ . If  $\gamma_1, \gamma_2, \dots, \gamma_j$  denote the distinct roots of the characteristic polynomial

$$u(x) = x^k - a_1 x^{k-1} - a_2 x^{k-2} - \dots - a_k$$

of the sequence  $U$ , then the terms of  $U$  can be written in the form

$$(1) \quad U_n = P_1(n) \gamma_1^n + P_2(n) \gamma_2^n + \dots + P_j(n) \gamma_j^n,$$

where  $P_i(n)$  ( $i=1, 2, \dots, j$ ) is a polynomial of  $n$  of degree less than the multiplicity of  $\gamma_i$  in  $u(x)$  and its coefficients are algebraic numbers from the field  $\mathbb{Q}(\gamma_1, \gamma_2, \dots, \gamma_j)$ .

Let  $G$  and  $H$  be two linear recurrences with characteristic polynomials

$$g(x) = x^k - A_1 x^{k-1} - A_2 x^{k-2} - \dots - A_k$$

and

$$h(x) = x^r - B_1 x^{r-1} - B_2 x^{r-2} - \dots - B_r,$$

respectively. P. Erdős asked (oral communication) whether the inequality

$$(2) \quad ||G_x| - |H_y|| < d$$

has finitely or infinitely many  $x, y$  solutions if  $d$  is a fixed positive number,  $G$  and  $H$  are not equivalent sequences and one of them is not a subsequence of the other.

In this paper we answer this question proving a general result from which it follows that under certain restrictions inequality (2) has only finitely many integer  $x, y$  solutions.

Let the distinct roots of the polynomials  $g(x)$  and  $h(x)$  be  $\alpha = \alpha_1, \alpha_2, \dots, \alpha_u$  and  $\beta = \beta_1, \beta_2, \dots, \beta_v$ , respectively. Throughout this paper we suppose that  $|\alpha| > |\alpha_2| \cong |\alpha_3| \cong \dots \cong |\alpha_u|$  and  $|\beta| > |\beta_2| \cong |\beta_3| \cong \dots \cong |\beta_v|$ , furthermore both  $\alpha$  and  $\beta$  have multiplicity one. Then  $|\alpha|, |\beta| > 1$  and by (1) we have

$$(3) \quad G_n = \alpha \alpha^n + P_2(n) \alpha_2^n + \dots + P_u(n) \alpha_u^n$$

and

$$(4) \quad H_n = \beta \beta^n + F_2(n) \beta_2^n + \dots + F_v(n) \beta_v^n,$$

1980 *Mathematics Subject Classification*. Primary 10A35; Secondary 10B45.

*Key words and phrases*. Recurrence, sequence of integers, differences, powers in sequences.

where  $a, b$  and the coefficients of the polynomials  $P_i(n)$  and  $F_i(n)$  are algebraic numbers. Furthermore let  $p_1 < p_2 < \dots < p_t$  be rational primes and denote by  $S$  the set of non-zero integers which have only these primes as prime factors with  $\pm 1 \in S$ . We prove:

**THEOREM 1.** *Suppose that  $G_i \neq a\alpha^i$ ,  $H_j \neq b\beta^j$ ,  $ab \neq 0$  and  $|s_1 a\alpha^i| \neq |s_2 b\beta^j|$  for  $i, j > n_0$  and for any integers  $s_1, s_2 \in S$ . Then*

$$(5) \quad ||s_1 G_x| - |s_2 H_y|| > \exp \{c \cdot \max(x, y)\}$$

for all integers  $x, y > n_1$  and  $s_1, s_2 \in S$ , where  $c$  and  $n_1$  are effectively computable positive numbers depending only on  $S$ ,  $n_0$  and the parameters of  $G$  and  $H$ .

This theorem implies several consequences, among others the answer to the question of Erdős.

**COROLLARY 1.** *Let  $G$  and  $H$  be linear recurrences satisfying the conditions of Theorem 1 with  $s_1 = s_2 = 1$  and let  $d$  be a positive number. Then the inequality*

$$||G_x| - |H_y|| < d$$

has only finitely many positive integer  $x, y$  solutions.

**COROLLARY 2.** *Let  $G$  be a linear recurrence satisfying the conditions  $a \neq 0$ ,  $G_i \neq a\alpha^i$  and  $a\alpha^i \notin S$  for  $i > n_0$ . Then  $|G_x - s| > e^{c'x}$  for all integers  $s \in S$  and  $x > n_1$ , where  $c'$  and  $n_1$  are effectively computable positive constants depending only on  $n_0$ ,  $S$  and  $G$ .*

Similar results were obtained by some other authors for second order sequences. K. Győry [6] and K. Győry, P. Kiss and A. Schinzel [7] showed that for Lucas and Lehmer sequences  $U_n \in S$  ( $n > 6$ ) holds only for finitely many terms of finitely many sequences. F. Beukers [4] and F. Beukers and R. Tijdeman [5] gave upper bounds for the number of solutions of the equation  $U_x = d$ , where  $U$  is a non degenerate binary recurrence of rational integers or complex numbers and  $d$  is fixed. T. N. Shorey [10] proved a closer result. He proved that (5) holds with fixed  $s_1$  and  $s_2$  for sequence  $G$  and  $H$  defined by  $G_n = a_1 \lambda^n + a_2 \mu^n$  and  $H_n = a_3 \lambda^n + a_4 \mu^n$ , where  $\lambda, \mu, a_1, a_2, a_3$  and  $a_4$  are algebraic numbers and  $\lambda/\mu$  is not a root of unity.

For general linear recurrences we proved in [8] that, under the assumptions of Theorem 1, the equation  $s_1 G_x - s_2 H_y = 0$  has only finitely many  $x, y$  solutions. Furthermore,  $G_x \in S$  holds only for finitely many integer  $x$ . We note that this second result follows also from a result of C. L. Stewart [12] who proved that the greatest prime divisor of  $G_x$  is greater than  $(1-\varepsilon) \log x$  if  $0 < \varepsilon < 1$  and  $x$  is sufficiently large. Recently M. Mignotte, T. N. Shorey and R. Tijdeman [9] have dealt with similar problems. Their results concern linear recurrence  $U$  of algebraic numbers and among others they proved: Let  $\gamma_1, \gamma_2, \dots, \gamma_j$  be the distinct roots of the polynomial  $u(x)$  with  $|\gamma_1| = |\gamma_2| = \dots = |\gamma_r| > |\gamma_{r+1}| \geq \dots \geq |\gamma_j|$ . Assume  $r \leq 3$ ,  $|\gamma_1| > 1$  and at least one of the numbers  $\gamma_i/\gamma_k$  with  $1 \leq i < k \leq r$  is not a root of unity. Then there exist computable numbers  $c > 0$  and  $c' > 0$  depending only on the sequence  $U$  such that

$$|U_m - U_n| \geq |\gamma_1|^m \exp \{-c (\log m)^2 \log(n+2)\}$$

whenever  $m > c'$  and  $m > n$ .

Theorem 1 and the corollaries give certain improvements and generalizations of the results mentioned above.

Another result for general linear recurrences was obtained by T. N. Shorey and C. L. Stewart [11]. They proved that if  $G$  satisfies the conditions of our Corollary 2 and  $d$  is a fixed non-zero integer, furthermore if

$$dx^q = G_n$$

for integers  $x$  and  $q$  larger than one, then  $q < N$ , where  $N$  is an effectively computable number depending only on  $d$  and  $G$ . We give two generalizations of this result.

**THEOREM 2.** *Let  $G$  be a linear recurrence satisfying the conditions  $a \neq 0$  and  $G_i \neq \alpha \alpha^i$  for  $i > n_0$ . If*

$$sw^y = G_x$$

*for integers  $s \in S, x, y$  and  $w$ , then  $y < N$ , where  $N$  is effectively computable in terms of  $n_0, S$  and the parameters of the sequence  $G$ .*

**THEOREM 3.** *Let  $G$  be a linear recurrence with conditions  $a \neq 0, |\alpha_2| \neq 1, |\alpha_2| > |\alpha_3|, P_2(i) \neq 0$  and  $G_i \neq \alpha \alpha^i$  for  $i > n_0$ . Then*

$$|sw^y - G_x| > e^{cx}$$

*for all integers  $s, w, x$  and  $y$  with  $s \in S$  and  $x, y > n_1$ , where  $c$  and  $n_1$  are effectively computable positive numbers depending only on  $n_0, S$  and the parameters of  $G$ .*

We denote that Theorem 3 does not hold rejecting the condition  $|\alpha_2| \neq 1$ . For example let us consider the second order recurrence  $G$  defined by  $G_n = 2^n - 1$ . In this case we have  $|sw^y - G_x| = 1$  with  $s = 1$  and  $w = 2$  for infinitely many integers  $x$  and  $y$ .

An easy consequence of Theorem 3 is as it follows.

**COROLLARY 3.** *Let  $G$  be a linear recurrence satisfying the conditions of Theorem 3 and let  $d$  be a non-negative integer. If*

$$|sw^y - G_x| = d$$

*for integers  $x, y, w$  and  $s \in S$ , then  $y < N$ , where  $N$  is effectively computable in terms of  $n_0, S, G$ , and  $d$ .*

For the proofs of our results we need two results due to A. Baker.

**THEOREM A.** *Let*

$$A = \gamma_0 + \gamma_1 \log \omega_1 + \gamma_2 \log \omega_2 + \dots + \gamma_n \log \omega_n,$$

*where the  $\gamma$ 's and the  $\omega$ 's denote algebraic numbers ( $\omega_i \neq 0$  or  $1$ ). We assume that not all the  $\gamma_i$ 's are zero and that the logarithms mean their principal values. Suppose that  $\omega_1$  and  $\gamma_i$  have heights at most  $M_i$  ( $\geq 4$ ) and  $B$  ( $\geq 4$ ), respectively, and that the field generated by the  $\omega$ 's and  $\gamma$ 's over the rational numbers has degree at most  $d$ . If  $A \neq 0$ , then*

$$|A| > (B\Omega)^{-C\Omega \log \Omega},$$

where

$$\Omega = \log M_1 \log M_2 \dots \log M_n,$$

$$\Omega' = \Omega / \log M_n$$

and  $C = (16nd)^{200n}$ . If  $\gamma_0 = 0$  and  $\gamma_1, \gamma_2, \dots, \gamma_n$  are rational integers, then

$$|A| > B^{-C\Omega \log \Omega'}.$$

(See A. Baker [2] or A. Baker and C. L. Stewart [3].)

**THEOREM B.** Let  $\omega_1, \omega_2, \dots, \omega_n$  be algebraic numbers as in Theorem A and let  $\gamma_0 = 0$ . Then there exists an effectively computable number  $C > 0$  depending only on  $n, d, M_1, \dots, M_{n-2}$  and  $M_{n-1}$  such that, for  $A$  defined in Theorem A and for any  $\delta$  with  $0 < \delta < 1/2$ , the inequalities

$$0 < |A| < (\delta/B)^{C \cdot \log M} e^{-\delta B}$$

have no solution in rational integers  $\gamma_1, \gamma_2, \dots, \gamma_{n-1}$  and  $\gamma_n (\neq 0)$  with absolute values at most  $B$  and  $B'$ , respectively (see A. Baker [1]).

Now we shall prove our results. In the proofs we shall denote by  $c_1, c_2, \dots, n_1, n_2, \dots$  positive numbers which are effectively computable in terms of the given parameters.

**PROOF OF THEOREM 1.** Let  $G$  and  $H$  be sequences satisfying the conditions. By (3) we have

$$G_x = a\alpha^x \left( 1 + \sum_{i=2}^n \frac{P_i(x)}{a} \left( \frac{\alpha_i}{\alpha} \right)^x \right)$$

therefore

$$(6) \quad e^{c_1 x} < |G_x| < e^{c_2 x}$$

and similarly

$$(7) \quad e^{c_3 y} < |H_y| < e^{c_4 y}$$

for  $x, y$  sufficiently large. Let  $c_5$  be a real number with  $0 < c_5 < \min(c_1, c_3, \log 2)$  and suppose that

$$(8) \quad ||s_1 G_x| - |s_2 H_y|| < e^{c_5 \max(x, y)}$$

for some integers  $s_1, s_2 \in S$  and  $x, y > n_0$ . Without loss of generality we can assume that  $(s_1, s_2) = 1$  and  $s_1 G_x$  and  $s_2 H_y$  have the same sign, say positive. Let  $m = \min(x, y)$ ,  $M = \max(x, y)$  and  $k = \max(k_1, k_2, \dots, k_t)$ , where the  $k_i$ 's are defined by  $s_1 s_2 = p_1^{k_1} p_2^{k_2} \dots p_t^{k_t}$ .

First we prove that  $k < M^3$  if  $M$  is large enough. We can suppose that  $k > M$  and  $s_1$  has the greatest prime factor with power  $k$ . It implies, by (6) and (8), that

$$(9) \quad \left| 1 - \frac{s_2 H_y}{s_1 G_x} \right| < \frac{e^{c_5 M}}{s_1 G_x} < e^{c_5 M - c_1 x - c_6 k} < e^{-c_7 M - c_1 x}$$

since  $c_6 \geq \log 2$  and so  $c_5 < c_6$  by the choice of  $c_5$ . It shows that  $|(s_2 H_y)/(s_1 G_x)|$

tends to 1 if  $M \rightarrow \infty$ , therefore

$$(10) \quad \left| \log \frac{s_2 H_y}{s_1 G_x} \right| < 2e^{c_5 M - c_1 x - c_6 k}$$

for  $M > n_3$ . Using Theorem A we give a lower estimation for  $|\log(s_2 H_y)/(s_1 G_x)|$ . In our case, using the notations of Theorem A,  $n \leq t+2$ ,  $\omega_1, \omega_2, \dots, \omega_{n-2}$  are primes ( $\in S$ ),  $\omega_{n-1} = H_y$ ,  $\omega_n = G_x$ ,  $\gamma_0 = 0$ ,  $B = k$ ,  $M_i = \max(p_1, \dots, p_t, 4)$  for  $1 \leq i \leq n-2$ ,  $M_{n-1} = H_y$ ,  $M_n = G_x$ ,  $\Omega = c_8 xy$  and  $\Omega' = c_9 M$ . The above mentioned result of [8] implies that  $(s_2 H_y)/(s_1 G_x) \neq 1$  for  $M > n_3$  so by Theorem A we have

$$(11) \quad \left| \log \frac{s_2 H_y}{s_1 G_x} \right| > k^{-c_{10} xy \log M} = e^{-c_{10} xy \log M \log k}.$$

It follows from (10) and (11) that

$$c_5 M - c_1 x - c_6 k + \log 2 > -c_{10} xy \log M \log k$$

and so

$$(12) \quad c_5 < c_1 \frac{M}{k} - c_6 \frac{x}{k} + \frac{1}{k} c_{10} xy \log M \log k + \frac{\log 2}{k}.$$

If  $k \geq M^3$  then  $\sqrt[3]{k} \geq M$  and  $k = \sqrt[3]{k}^2 \sqrt[6]{k}^2 \geq M^2 \sqrt[6]{M} \sqrt[6]{k}$  which imply the inequalities  $c_5(M/k) < c_6/4$ ,  $c_1(x/k) < c_6/4$ ,  $(\log 2)/k < c_6/4$  and

$$\frac{1}{k} c_{10} xy \log M \log k \leq c_{10} \frac{x}{M} \frac{y}{M} \frac{\log M}{\sqrt[6]{M}} \frac{\log k}{\sqrt[6]{k}} < c_6/4$$

for  $M > n_4$ . But these inequalities contradict to (12) which proves the assertion. In what follows we assume  $M = x \geq y$ . Choosing  $c_{11} = c_1 + c_7$ , by (9) we have

$$(13) \quad 1 - e^{-c_{11} x} < \frac{s_2 H_y}{s_1 G_x} < 1 + e^{-c_{11} x}.$$

But using the explicit form of  $G_x$  we get

$$1 - e^{-c_{11} x} < |G_x/(a\alpha^x)| < 1 + e^{-c_{11} x}$$

for  $x > n_5$  since  $|x_i/\alpha| < 1$  ( $i = 2, 3, \dots, u$ ). From this by (13) it follows that

$$(14) \quad (1 - e^{-c_{11} x})(1 - e^{-c_{13} x}) < \left| \frac{s_2 H_y}{s_1 a \alpha^x} \right| < (1 + e^{-c_{11} x})(1 + e^{-c_{13} x})$$

and so, on taking logarithms, we obtain

$$(15) \quad \left| \log \left| \frac{s_2 H_y}{s_1 a \alpha^x} \right| \right| < e^{-c_{13} x}$$

for  $x$  sufficiently large. Here we can assume that

$$\left| \frac{s_2 H_y}{s_1 a \alpha^x} \right| \neq 1.$$

For if  $|s_2 H_y| = |s_1 a \alpha^x|$  then  $s_1 a \alpha^x$  is an algebraic integer and it is divisible by each prime factor of  $H_y$ . But, as it was mentioned above, Stewart [12] proved under our conditions that  $H_y$  has a prime factor larger than  $(1/2) \log y$  for  $y$  large enough. It is a contradiction for  $y > n_6$  since  $s_1 a \alpha^x$  can be divisible only by finitely many primes. Using again Theorem A with  $B = x^3$ , the following estimation can be given:

$$(16) \quad \left| \log \left| \frac{s_2 H_y}{s_1 a \alpha^x} \right| \right| > e^{-c_{14} y \log x}.$$

Comparing (15) with (16) we get

$$(17) \quad y > c_{15} \frac{x}{\log x}.$$

Now we shall give, similarly as above, a lower and an upper estimation for the absolute value of the logarithm of

$$E = \left| \frac{s_2 b \beta^y}{s_1 a \alpha^x} \right|,$$

where  $E \neq 1$  by the conditions. First we know by (14) that

$$(18) \quad \frac{(1 - e^{-c_{11}x})(1 - e^{-c_{12}x})}{1 + e^{-c_{16}y}} < E < \frac{(1 + e^{-c_{11}x})(1 + e^{-c_{12}x})}{1 - e^{-c_{16}y}}$$

since

$$1 - e^{-c_{16}y} < \left| \frac{H_y}{b \beta^y} \right| = \left| 1 + \sum_{i=2}^v \frac{F_i(y)}{b} \left( \frac{\beta_i}{\beta} \right)^y \right| < 1 + e^{-c_{16}y}$$

if  $x$  and so by (17)  $y$  is large. We have assumed that  $y \leq x$ , therefore by (18) we have

$$|\log E| < e^{-c_{17}y}.$$

On the other hand by Theorem A we get

$$|\log E| > e^{-c_{18} \log x}.$$

From the last two inequalities

$$y < c_{19} \log x$$

follows which contradicts to (17) if  $x > n_7$ . Thus inequality (8) does not hold for  $x, y$  sufficiently large which proves the theorem with  $c = c_5$ .

**PROOF OF COROLLARY 1.** By Theorem 1 we have to prove that the inequality

$$|G_x| < d + |H_y|$$

does not hold for  $x$  sufficiently large and  $y \leq n_1$ . But it is obvious by (6).

PROOF OF COROLLARY 2. It can be proved similarly as Theorem 1 using that  $G_x \notin S$  for  $x$  sufficiently large, which was proved in [8], furthermore that  $s/(a\alpha^x) \neq 1$  in view of the conditions.

The proof of Theorem 2 follows from the proof of Theorem 3 therefore now we prove first Theorem 3.

PROOF OF THEOREM 3. Let  $0 < c < \min(c', c_5, (\log |\alpha_2|)/2)$ , where  $c'$  and  $c_5$  are defined in Corollary 2 and in the proof of Theorem 1, respectively. Further, let  $s \in S$ ,  $x > n_0$ ,  $y$  and  $w$  be rational integers such that

$$(19) \quad |sw^y - G_x| \leq e^{cx}.$$

We may assume that  $s, w^y$  and  $G_x$  are positive integers. Furthermore we assume that  $p_i^r \nmid s$  for  $r \geq y$  and  $i = 1, 2, \dots, t$ , and so  $s < (p_1 p_2 \dots p_t)^y$ , since otherwise the perfect  $y$ -th powers of  $p_i^r$ 's could increase the value of  $w^y$ . Then by (19) and (6) we have

$$(20) \quad sw^y < G_x + e^{cx} < e^{c_{20}x}$$

for  $x > n_8$ . It implies the inequality  $y < c_{21}x$  since  $w > 1$  by Corollary 2. From (19) it follows that

$$\left| \frac{sw^y}{G_x} - 1 \right| \leq \frac{e^{cx}}{G_x} \leq e^{-c_{22}x}$$

and we get in the same way as in the proof of Theorem 1 that

$$(21) \quad (1 - e^{-c_{12}x})(1 - e^{-c_{22}x}) < \left| \frac{sw^y}{a\alpha^x} \right| < (1 + e^{-c_{12}x})(1 + e^{-c_{22}x}).$$

We show that  $|(sw^y)/(a\alpha^x)| \neq 1$ . Suppose the converse that is  $sw^y = |a\alpha^x|$ . One can easily see that, under our assumptions,  $a\alpha^x$  is a real number and that (19) does not hold for large  $x$  in case  $sw^y = -a\alpha^x$ . Thus we may assume that  $sw^y = a\alpha^x$  and in this case for any  $\varepsilon > 0$  and  $x$  sufficiently large we have

$$|sw^y - G_x| = |P_2(x)\alpha_2^x| \left| 1 + \sum_{i=3}^n \frac{P_i(x)}{P_2(x)} \left( \frac{\alpha_i}{\alpha_2} \right)^x \right| \geq e^{(\log |\alpha_2| - \varepsilon)x}$$

which contradicts to (19) if  $|\alpha_2| > 1$ . In case  $|\alpha_2| < 1$  we also get a contradiction since  $|sw^y - G_x|$  is an integer, but the last equation implies that  $0 < |sw^y - G_x| < 1$ .

Thus  $|(sw^y)/(a\alpha^x)| \neq 1$  for large  $x$  and by (21)

$$(22) \quad 0 < \left| \log \left| \frac{sw^y}{a\alpha^x} \right| \right| < e^{-c_{23}x}.$$

Now we apply Theorem B with  $\omega_n = w$ ,  $M_n = w$ ,  $B = c_{24}x$  (since the exponents of  $p_i$ 's in  $s$  are less than  $c_{24}x$  by (20)) and  $B' = y$ . We get

$$(23) \quad \left| \log \left| \frac{sw^y}{a\alpha^x} \right| \right| > e^{-c_{25} \cdot (\log y - \log \delta) \cdot \log w - c_{24}\delta x}.$$



From (22) and (23) it follows that

$$c_{25}(\log y - \log \delta) \log w + c_{24} \delta x > c_{23} x$$

that is with  $0 < \delta < c_{23}/c_{24}$

$$(24) \quad c_{25}(\log y - \log \delta) \log w > c_{26} x.$$

But by (20)

$$y \log w < c_{20} x$$

follows which, together with (24), implies the inequality

$$\log y - \log \delta > c_{27} y.$$

It holds only for finitely many integer  $y$ .

Thus for  $x$  sufficiently large (19) implies  $y < n_0$  which proves the assertion of Theorem 3.

PROOF OF THEOREM 2. Let

$$(25) \quad sw^y = G_x$$

for some integers  $x > n_0$ ,  $y$ ,  $w$  and  $s \in S$ . In [8] we have proved that  $G_x \in S$  holds only for finitely many integer  $x$ , therefore  $w > 1$ . As in the proof of Theorem 3 we have

$$\left| \frac{sw^y}{ax^x} - 1 \right| = \left| \sum_{i=2}^u \frac{P_i(x)}{a} \left( \frac{\alpha_i}{\alpha} \right)^x \right| < e^{-c_{28}x}.$$

Here  $|(sw^y)/(ax^x)| \neq 1$  since otherwise by (25)  $G_x = ax^x$  would follow which contradicts to the conditions.

In the sequel the proof can be carried out similarly as in the proof of Theorem 3.

## REFERENCES

- [1] BAKER, A., A sharpening of the bounds for linear forms in logarithms II, *Acta Arithm.* **24** (1973), 33—36. *MR* **51**#12724.
- [2] BAKER, A., The theory of linear forms in logarithms, *Transcendence theory: Advances and Applications*, Acad. Press, London and New York, 1977. *MR* **58**#16543.
- [3] BAKER, A. and STEWART, C. L., Further aspects of transcendence theory, *Astérisque* **41—42** (1977), 153—163. *MR* **56**#5444.
- [4] BEUKERS, F., The multiplicity of binary recurrences, *Compositio Math.* **40** (1980), 251—267. *MR* **81g**: 10019.
- [5] BEUKERS, F. and TIJDEMAN, R., On the multiplicities of binary complex recurrences, *Compositio Math.* **51**(1984), 193—213. *MR* **85i**: 11017.
- [6] GYÖRY, K., On some arithmetical properties of Lucas and Lehmer numbers, *Acta Arithm.* **40** (1982), 369—373. *MR* **83k**: 10018.
- [7] GYÖRY, K., KISS, P. and SCHINZEL, A., A note on Lucas and Lehmer sequences and their applications to Diophantine equations, *Colloq. Math.* **45** (1981), 75—80. *MR* **83g**: 10009.
- [8] KISS, P., On common terms of linear recurrences, *Acta Math. Acad. Sci. Hungar.* **40** (1982), 119—123. *MR* **84h**:10014.

- [9] MIGNOTTE, M., SHOREY, T. N. and TIJDEMAN, R., The distance between terms of an algebraic recurrence sequence, *J. Reine Angew. Math.* **349**(1984), 63—76. *MR 85f*: 11007.
- [10] SHOREY, T. N., Linear forms in members of a binary recursive sequence, *Acta Arith.* **43**(1984), 317—331. *MR 85m*: 11038.
- [11] SHOREY, T. N. and STEWART, C. L., On the Diophantine equation  $ax^{2t} + bx^t y + cy^2 = d$  and pure powers in recurrence sequences, *Math. Scand.* **52** (1983), 24—36. *MR 84g*: 10038.
- [12] STEWART, C. L., On divisors of terms of linear recurrence sequences, *J. Reine Angew. Math.* **333** (1982), 12—31. *MR 83i*: 10010.

(Received August 12, 1983)

TANÁRKÉPZŐ FŐISKOLA  
MATEMATIKAI TANSZÉK  
LEÁNYKA U. 4.  
H-3300 EGER  
HUNGARY



# A NOTE ON THE MÖBIUS AND LIOUVILLE FUNCTIONS

G. HARMAN, J. PINTZ and D. WOLKE

## 1. Introduction

The purpose of this short note is to investigate the values taken by the Möbius function  $\mu(n)$  on consecutive square-free integers, and by the Liouville function  $\lambda(n)$  on successive integers. We recall the definition of these two functions:

$\lambda(1)=1$ ;  $\lambda(n)=(-1)^{\Omega(n)}$ , where  $\Omega(n)$  denotes the number of prime factors of  $n$  counted according to multiplicity.

$\mu(n)=\lambda(n)$  if  $n$  is square-free (henceforth abbreviated to sf).  
 $=0$  if  $d^2|n$  with  $d>1$ .

We write  $q_k$  for the  $k$ -th sf integer. Also we use  $\sigma$  and  $s$  to denote  $+$  or  $-$  signs. With this notation we define

$$(1) \quad f_{\sigma s}(x) = \sum_{\substack{k \leq x \\ \mu(q_k) = \sigma 1 \\ \mu(q_{k+1}) = s 1}} 1; \quad g_{\sigma s}(x) = \sum_{\substack{n \leq x \\ \lambda(n) = \sigma 1 \\ \lambda(n+1) = s 1}} 1.$$

Since we have no reason *a priori* to assume the contrary, we may plausibly state the following

CONJECTURE. For each case  $(s, \sigma) = (+, +), (+, -), (-, -), (-, +)$  we have

$$(2) \quad f_{\sigma s} \sim \frac{x}{4} \quad \text{and} \quad g_{\sigma s} \sim \frac{x}{4} \quad \text{as} \quad x \rightarrow \infty.$$

We further write  $f^*(x) = f_{++}(x) + f_{--}(x)$  and define  $g^*(x)$  analogously. Also we put

$$M(X) = \sum_{n \leq X} \mu(n), \quad L(X) = \sum_{n \leq X} \lambda(n).$$

It is then easy to verify the following simple relations between the functions defined in (1):

$$(3) \quad f_{++}(x) = f_{--}(x) + M(q_x) + \frac{1}{2}(\mu(q_{x+1}) - 1)$$

$$f_{+-}(x) = f_{-+}(x) + \frac{1}{2}(1 - \mu(q_{x+1}))$$

and similarly for  $g_{\sigma s}$  *mutatis mutandis*. It is well-known that

$$(4) \quad M(X) = o(X), \quad L(X) = o(X), \quad q_x = \pi^2 x/6 + O(x^{1/2}),$$

(Theorems 333 and 335 of [1] give the first and last of the results in (4) for example). Hence by (3) and (4) it suffices to show that for one pair  $(\sigma, s)$  we have  $f_{\sigma s}(x) \sim x/4$  in order to prove our conjecture for all four  $f$  functions and similarly for  $g_{\sigma s}$ . The conjecture (2) is clearly equivalent to the supposition that

$$\sum_{k \leq x} \mu(q_k) \mu(q_{k+1}) = o(x), \quad \sum_{n \leq x} \lambda(n) \lambda(n+1) = o(x).$$

To see this consider, for example,  $g_{+-}(x)$  for which we have

$$g_{+-}(x) = \frac{1}{4} \sum_{n \leq x} (\lambda(n)+1)(1-\lambda(n+1)) = \frac{1}{4} ([x] - \sum_{n \leq x} \lambda(n) \lambda(n+1) + 1 - \lambda([x]+1)).$$

These conjectures are apparently in the same league of difficulty as the prime twins problem so we are very far from proving them! In the rest of the paper, however, we will demonstrate that it is possible to obtain non-trivial bounds for our functions. Our two results are as follows:

**THEOREM 1.** *For the functions defined by (1), and any  $\varepsilon > 0$ , we have*

$$(5) \quad f_{\sigma s}(x) \equiv x(\log x)^{-7-\varepsilon}; \quad g_{\sigma s}(x) \equiv x(\log x)^{-7-\varepsilon}$$

for  $(\sigma, s) = (+, -)$  and  $(-, +)$  and all sufficiently large  $x$ .

**THEOREM 2.** *For the functions defined by (1) we have, for  $(\sigma, s) = (+, +)$  and  $(-, -)$ ,*

$$(6) \quad f_{\sigma s}(x) > \frac{x}{60} (1 + o(1)); \quad g_{\sigma s} \equiv \frac{x}{6} (1 + o(1)).$$

## 2. Proof of Theorem 1

First we mention a completely elementary argument. By Theorem 418 of [1] for every  $n \geq 1$  the interval  $(n, 2n]$  contains a prime, and so every interval of the form  $(2n, 4n]$  for  $n \geq 2$  contains a prime and a square-free integer which is twice a prime. It follows that

$$f_{\sigma s}(x) \gg \log x, \quad g_{\sigma s}(x) \gg \log x.$$

The somewhat surprising feature of this result is that our proof of Theorem 1 is just as unsophisticated!

In 1979 D. Wolke [6], improving work of D. R. Heath-Brown [3] and Y. Motohashi [4] established, by analytic methods, that almost all intervals of the form

$$[n, n + (\log n)^C)$$

contain  $sf$  numbers with two odd prime factors. He gave the value  $5.10^6$  for  $C$ . The "almost all" in the result indicates that the number of integers  $\leq x$  for which the given interval does not contain almost-primes is  $o(X)$ . More recently, G. Harman

[2] improved the value of  $C$  to  $7+\eta$  for any  $\eta>0$ . The important point for our present investigation is that, as opposed to a sieve method [3], the numbers obtained do have precisely two prime factors. It follows that almost all intervals of the form  $[n, n+2(\log n)^C]$  contain sf integers of both forms  $p_1 p_2, 2p_3 p_4$  ( $p_j$  primes). Hence there must be a  $q_x$  in the interval with  $\mu(q_x)\mu(q_{x+1})=-1$ . By the relations (3) the result follows for  $f_{+-}$  and  $f_{-+}$  taking  $\eta=\varepsilon/2$  and  $x$  sufficiently large. The result follows *a fortiori* for  $g_{+-}$  and  $g_{-+}$ . It should be noted that one expects there to be many  $n$  in the intervals considered with  $\lambda(n)\lambda(n+1)=-1$ , and perhaps it is possible to construct an argument that should there only be a few such  $n$  in certain intervals then there must be many in other intervals (cf. our proof of Theorem 2 below).

### 3. Proof of Theorem 2

First we note that analytic methods give results (see [5]) which only imply lower bounds of the form  $\sqrt{x} \exp(-3(\log \log x)^{3/2})$  for  $f_{++}, f_{--}$  and  $g_{--}$ , and no result at all for  $g_{++}$ . However, it was this observation together with the argument of Section 2 which prompted this investigation.

In this section our arguments will again be very elementary. We would like to express our thanks to Dr. R. R. Hall for showing us a simple proof that

$$\sum_{n \leq X} \lambda(n)\lambda(n+1) \cong -\frac{2X}{3} + O(1).$$

This was produced as an answer to a question put by E. Bombieri at the Durham Conference on Analytic Number Theory in 1979. Our proof given here was inspired by Hall's proof and leads to the stronger bound

$$(7) \quad \sum_{n \leq X} \lambda(n)\lambda(n+1) \cong -\frac{X}{3} + O(\log X).$$

We shall in fact prove that

$$(8) \quad f^*(x) > \frac{x}{30} + O(\sqrt{x}),$$

and

$$(9) \quad g^*(x) \cong \frac{x}{3} + O(\log x).$$

We note (6) follows from (3), (4), (8) and (9), while (7) follows immediately from (9) (recalling the definition of  $f^*$  and  $g^*$  in §1).

We start by proving (9). Let us suppose that  $\lambda(x)\lambda(x+1)=-1$ . Hence  $\lambda(2x)\lambda(2x+2)=-1$  so either  $\lambda(2x)\lambda(2x+1)=1$  or  $\lambda(2x+1)\lambda(2x+2)=1$ . By considering all integers  $x \in [X, 2X]$  it follows that

$$g^*(4X) - g^*(X) \cong X + O(1),$$

and so, by induction,

$$g^*(x) \cong \frac{x}{3} + O(\log x).$$

To use a similar argument for  $f^*(x)$  we need to be more careful. The following lemma is a particular example of a well-known result: — *the sf equivalent to the prime  $k$ -tuple conjecture is true*. We give the proof (cf. the proof of Theorem 333 in [1]) for completeness. While writing this paper we learnt that D. R. Heath-Brown was aware of a similar argument.

LEMMA. Write

$$P_1(y) = (36y+3)(36y+5)(72y+7),$$

$$P_2(y) = (36y+31)(36y+33)(72y+65).$$

Then the number of sf integers of the form  $P_1(y)$  or  $P_2(y)$  for  $1 \leq y \leq N$  is

$$\frac{2N(27-2\pi^2)}{\pi^2} + O(\sqrt{N}).$$

PROOF. We consider integers of the form  $P_1(y)$ , the proof for  $P_2(y)$  follows analogously and, of course,  $P_1(y_1) = P_2(y_2)$  has no solutions in integers.

We write

$$A(N) = \sum_{y \leq N} |\mu(P_1(y))|.$$

Then

$$(10) \quad A(N) = \sum_{1 \leq y \leq N} \sum_{d^2 | P_1(y)} \mu(d) = \sum_{d < 9N^{1/2}} \mu(d) \sum_{\substack{y=1 \\ d^2 | P_1(y)}}^N 1.$$

If we let  $Y(N, d)$  denote the sum over  $y$  on the right of (10) we then have  $Y(N, 1) = N$ ;  $Y(N, d) = 0$  if  $3|d$  or  $2|d$ . On the other hand, by considering congruences modulo  $d^2$ ,  $Y(N, d) = 3N/d^2 + O(1)$  if  $(6, d) = 1$  and  $d > 1$  (note the three factors of  $P_1(y)$  are always coprime). Hence

$$\begin{aligned} A(N) &= 3N \sum_{\substack{(d, 6)=1 \\ 1 \leq d \leq 9N^{1/2}}} \frac{\mu(d)}{d^2} - 2N + O(\sqrt{N}) = \\ &= 3N \sum_{\substack{d=1 \\ (d, 6)=1}}^{\infty} \frac{\mu(d)}{d^2} - 2N + O(\sqrt{N}) = \\ &= \frac{3N}{\zeta(2)} \frac{4}{3} \frac{9}{8} - 2N + O(\sqrt{N}) = \\ &= \frac{N(27-2\pi^2)}{\pi^2} + O(\sqrt{N}). \end{aligned}$$

This completes the proof of the lemma.

It is now easy to establish (8) and thus complete the proof of Theorem 2. By our lemma the number of sf triplets of the form  $36n+3$ ,  $36n+5$ ,  $72n+7$  or  $36n+31$ ,  $36n+33$ ,  $72n+65$  less than  $q_x$  (recall (4) for the size of  $q_x$ ) is

$$\frac{(27-2\pi^2)}{216} x + O(\sqrt{x}) > \frac{x}{30} + O(\sqrt{x}).$$



We note that  $\mu(m)=0$  for  $m$  of the form  $36n+4$ ,  $72n+8$ ,  $72n+9$  etc. Now if  $\mu(36n+3)\mu(36n+5)=-1$  it follows that either  $\mu(72n+6)\mu(72n+7)=1$  or  $\mu(72n+7)\mu(72n+10)=1$ . A similar argument holds for the triplets of the other form. Thus

$$f^*(x) > \frac{x}{30} + O(\sqrt{x})$$

and the proof is finished. It is, of course, possible to improve the 60 in our result. For example one can consider, in addition to the above argument, sf quintuplets of the form  $72y+41$ ,  $72y+42$ ,  $72y+43$ ,  $144y+83$ ,  $144y+85$ , since corresponding to each quintuplet there exists a  $q_x$  (not of a type given by our previous discussion), with  $\mu(q_x)\mu(q_{x+1})=1$ .

## REFERENCES

- [1] HARDY, G. H. and WRIGHT, E. M., *An introduction to the theory of numbers*, Oxford, 5th edition, 1979. MR 81i: 10002.
- [2] HARMAN, G., Almost-primes in short intervals, *Math. Ann.* 258 (1981), 107—112. MR 83d: 10048.
- [3] HEATH-BROWN, D. R., Almost-primes in arithmetic progressions and short intervals, *Math. Proc. Cambridge Philos. Soc.* 83 (1978), 357—375. MR 58# 10789.
- [4] MOTOHASHI, Y., A note on almost-primes in short intervals, *Proc. Japan Acad. Ser. A Math. Sci.* 55 (1979), 225—226. MR 81a: 10056.
- [5] PINTZ, J., Oscillatory properties of  $M(x) = \sum_{n \leq x} \mu(n)$  III, *Acta Arith.*
- [6] WOLKE, D., Fast-Primzahlen in kurzen Intervallen, *Math. Ann.* 244 (1979), 233—242. MR 81i: 10055.

(Received September 5, 1983)

IMPERIAL COLLEGE  
LONDON SW7  
ENGLAND

MTA MATEMATIKAI KUTATÓ INTÉZET  
P.O. BOX 127  
H-1364 BUDAPEST  
HUNGARY

FACHBEREICH MATHEMATIK  
UNIVERSITÄT FREIBURG  
ALBERTSTRASSE 23B  
D-78 FREIBURG  
FEDERAL REPUBLIC OF GERMANY



# ON THE SUM OF A $\iota$ -SEMILATTICE ORDERED SYSTEM OF ALGEBRAS

JERZY PŁONKA

0. In [7] the notion of the sum of a direct system of algebras was introduced. This notion has been considered in some papers (see e.g. [2], [3], [4], [5], [6], [8], [9]). However, in [7] we studied only algebras without nullary polynomials. This inconvenience was removed in [10] where we assumed that the values of nullary polynomials belong to the component with the least index  $\omega$ . We proved in [10] that such sum preserves exactly all regular identities satisfied in any component. Recall that the identity  $\varphi = \psi$  is called regular if the sets of variables of  $\varphi$  and  $\psi$  are identical. However, we can ask what happens if we assume that the values of nullary polynomials belong to the component with the greatest index  $\iota$ . This we do in this paper, namely we define a construction called the sum of a  $\iota$ -semilattice ordered system of algebras.

A regular identity not containing nullary fundamental operation symbols will be called strictly regular and an identity with nullary fundamental operation symbols occurring on both sides will be called nullary symmetrical. In Section 1 we prove that in the case if this construction is not trivial it preserves all strictly regular identities satisfied in any component, preserves all nullary symmetrical identities satisfied in the component with the greatest index  $\iota$  and does not preserve any other (see Theorem 1).

In Section 2 we prove that: if  $K$  is a variety of type  $\tau$  such that an equality  $x \circ y = x$  is satisfied in  $K$  for some term  $x_1 \circ x_2$  of type  $\tau$  not containing nullary fundamental operation symbols then the variety  $K_J$  of type  $\tau$  defined by all strictly regular and nullary symmetrical equalities satisfied in  $K$  — contains exactly sums of  $\iota$ -semilattice ordered systems of algebras belonging to  $K$  or  $K_U$  where  $K_U$  is the variety defined by identities satisfied in  $K$  and not containing nullary polynomial symbols (see Theorem 3). We also prove that under this assumptions  $K_J$  covers  $K$ .

1. Let us fix a type  $\tau: T \rightarrow N$  of algebras where  $\tau(T) \setminus \{0, 1\} \neq \emptyset$ . Denote  $T^* = \{t: t \in T, \tau(t) \neq 0\}$ ,  $\tau^* = \tau|_{T^*}$ ,  $T_0 = T \setminus T^*$ . If  $\mathfrak{A} = (A; \{f_t\}_{t \in T})$  is an algebra of the type  $\tau$ , we denote by  $\mathfrak{A}^*$  the reduct  $(A; \{f_t\}_{t \in T^*})$  of  $\mathfrak{A}$ . So  $\mathfrak{A}^*$  is of the type  $\tau^*$ . Let  $T' \subseteq T$ ,  $\tau' = \tau|_{T'}$ . If  $\mathfrak{A}_1 = (A_1; \{f_t\}_{t \in T'})$  is an algebra of type  $\tau'$  and  $\mathfrak{A}_2 = (A_2; \{f_t\}_{t \in T})$  is an algebra of the type  $\tau$  then the mapping  $h: A_1 \rightarrow A_2$  will be called a  $T'$ -homomorphism of  $\mathfrak{A}_1$  into  $\mathfrak{A}_2$  if  $h$  is a homomorphism of  $\mathfrak{A}_1$  into the reduct  $(A_2; \{f_t\}_{t \in T'})$  of  $\mathfrak{A}_2$ .

1980 *Mathematics Subject Classification*. Primary 08A05; Secondary 08A30.

*Key words and phrases*. Semilattice ordered systems of algebras, regular identities, varieties of algebras.

By a  $\iota$ -semilattice ordered system of algebras we shall mean a triple  $\mathcal{A} = \langle (I; \leq), \{\mathfrak{A}_i\}_{i \in I}, \{h_i^j\}_{i, j \in I, i \leq j} \rangle$  satisfying the following conditions (i)–(v):

(i)  $I$  is a non empty set partially ordered by the relation  $\leq$  having the least upper bound property — (l.u.b.)

(ii) if  $T_0 \neq \emptyset$  then there exists in  $I$  the greatest element  $\iota$ .

(iii)  $\{\mathfrak{A}_i\}_{i \in I}$  is a family of algebras where for  $i \neq \iota$  the algebra  $\mathfrak{A}_i$  is an algebra of the type  $\tau^*$  and  $\mathfrak{A}_i = (A_i; \{f_i^t\}_{t \in T^*})$  whereas  $\mathfrak{A}_\iota = (A_\iota; \{f_i^t\}_{t \in T})$  is an algebra of type  $\tau$ .

(iv) if  $i, j \in I$ ,  $i \neq j$  then  $A_i \cap A_j = \emptyset$ .

(v) any  $h_i^j$  is a  $T^*$ -homomorphism of  $\mathfrak{A}_i$  into  $\mathfrak{A}_j$  and the set  $\{h_i^j\}_{i, j \in I, i \leq j}$  satisfies two conditions any  $h_i^i$  is the identity map,  $h_i^j \cdot h_j^k = h_i^k$  for any  $i \leq j \leq k$ .

We define a new algebra  $\mathfrak{S}(\mathcal{A})$  of type  $\tau$  putting  $\mathfrak{S}(\mathcal{A}) = (\bigcup_{i \in I} A_i; \{f_i^t\}_{t \in T})$

where the operations  $f_i^t$  are defined as follows:  $f_i^t = f_\iota^t$  if  $\tau(t) = 0$ . If  $\tau(t) \neq 0$ ,  $a_k \in A_{i_k}$  for  $k = 1, \dots, \tau(t)$  and  $q = \text{l.u.b.}(i_1, \dots, i_{\tau(t)})$  then

$$f_i^t(a_1, \dots, a_{\tau(t)}) = f_i^t(h_{i_1}^q(a_1), h_{i_2}^q(a_2), \dots, h_{i_{\tau(t)}}^q(a_{\tau(t)})).$$

The algebra  $\mathfrak{S}(\mathcal{A})$  will be called the sum of a  $\iota$ -semilattice ordered system  $\mathcal{A}$  of algebras. We have:

(vi) if we consider algebras without nullary polynomials then  $\mathfrak{S}(\mathcal{A})$  coincides with the sum of a direct system of algebras defined in [7].

(vii) for  $i \neq \iota$  any  $A_i$  is a subalgebra of  $(\mathfrak{S}(\mathcal{A}))^*$  and  $A_\iota$  is a subalgebra of  $\mathfrak{S}(\mathcal{A})$ .

REMARK 1. Obviously  $\mathcal{A}$  and  $\mathfrak{S}(\mathcal{A})$  can be also defined exactly in the same way, if we assume that any  $\mathfrak{A}_i$  is of the type  $\tau$  and  $h_i^j$  are  $T$ -homomorphism. However, one can easily observe that for  $i \neq \iota$  the nullary operations  $f_i^t$  have no influence on the definition of  $f_i^t$ . Moreover, such definition of  $\mathfrak{S}(\mathcal{A})$  would complicate the representation theorem (see Theorem 3, Section 2).

Now we want to answer the question what kind of identities preserves the construction  $\mathfrak{S}(\mathcal{A})$ . First we need some notions. If  $\mathfrak{A}$  is an algebra of type  $\tau$ , we denote by  $E(\mathfrak{A})$  the set of all identities of the variety generated by  $\mathfrak{A}$ . If  $\varphi$  is a term of type  $\tau$ ,  $\{x_{k_1}, \dots, x_{k_m}\}$  is the set of all variables occurring in  $\varphi$ ,  $\{c_{j_1}, \dots, c_{j_n}\}$  is the set of all nullary fundamental operation symbols in  $\varphi$  then we shall write  $\varphi = \varphi(x_{k_1}, \dots, x_{k_m}, c_{j_1}, \dots, c_{j_n})$ . Further we shall denote by  $\varphi^{\mathfrak{A}}$  the realization of a term  $\varphi$  in the algebra  $\mathfrak{A}$ .

LEMMA 1. If  $\mathcal{A} = \langle (I; \leq), \{\mathfrak{A}_i\}_{i \in I}, \{h_i^j\}_{i, j \in I, i \leq j} \rangle$ ,  $\varphi = \varphi(x_{k_1}, \dots, x_{k_m}, c_{j_1}, \dots, c_{j_n})$  is a term of type  $\tau$   $a_{k_r} \in A_{q_r}$  for  $r = 1, \dots, m$ ,  $q = \text{l.u.b.}(q_1, \dots, q_m)$ , then

$$\varphi(a_{k_1}, \dots, a_{k_m}, c_{j_1}, \dots, c_{j_n}) = \begin{cases} \varphi^q(h_{q_1}^q(a_{k_1}), \dots, h_{q_m}^q(a_{k_m})) & \text{if } \{c_{j_1}, \dots, c_{j_n}\} = \emptyset \\ \varphi^i(h_{q_1}^i(a_{k_1}), \dots, h_{q_m}^i(a_{k_m}), c_{j_1}^i, \dots, c_{j_n}^i) & \text{otherwise.} \end{cases}$$

Proof is by the standard induction on the complexity of  $\varphi$ .

THEOREM 1. If  $|I| > 1$  then the set  $E(\mathfrak{S}(\mathcal{A}))$  contains all strictly regular equalities satisfied in any  $\mathfrak{A}_i$  for  $i \in I$ , contains all nullary symmetrical equalities satisfied in  $\mathfrak{A}_\iota$  and does not contain any other.

The proof follows from Lemma 1 and is similar to that of Theorem 1 from [7]. So we do not present it here.

Let  $\mathfrak{A}=(A; \{f_i\}_{i \in T})$  be an algebra of type  $\tau$ . A binary function  $\circ : A^2 \rightarrow A$  will be called a  $\tau$ -partition function or briefly a  $\tau$ -p-function of  $\mathfrak{A}$  if it satisfies the following identities:

- (1)  $x \circ x = x,$
- (2)  $(x \circ y) \circ z = x \circ (y \circ z),$
- (3)  $x \circ y \circ z = x \circ z \circ y$
- (4)  $f_i(x_1, \dots, x_{\tau(i)}) \circ y = f_i(x \circ y, \dots, x_{\tau(i)} \circ y), \quad i \in T$
- (5)  $f_i(x_1, \dots, x_{\tau(i)}) \circ x_k = f_i(x_1, \dots, x_{\tau(i)}), \quad (k \in \{1, \dots, \tau(i)\}), \quad i \in T$
- (6)  $y \circ f_i(x_1, \dots, x_{\tau(i)}) = y \circ f_i(y \circ x_1, \dots, y \circ x_{\tau(i)}), \quad i \in T$
- (7)  $x \circ f_i(x, \dots, x) = x \quad \text{for } i \in T \setminus T_0.$

REMARK 2. These identities coincide with those of (1)–(7) from [7] if  $T_0 = \emptyset$ . So if  $T_0 = \emptyset$  then  $\tau$ -p-function coincides with p-function. Observe that if  $i \in T_0$  then (4) is non-regular.

Similarly like in [7] we can formulate the following

**THEOREM 2.** *To every  $\tau$ -p-function of the algebra  $\mathfrak{A}=(A; \{f_i\}_{i \in T})$  there corresponds a representation  $\mathfrak{A}=\mathfrak{S}(\mathcal{A})$  obtained as follows. Divide  $A$  into disjoint subsets  $A_i$  ( $i \in I$ ) putting two elements  $a, b$  of  $A$  into the same set  $A_i$  if and only if  $a \circ b = a$  and  $b \circ a = b$ . In the set  $I$  of indices we introduce the relation " $\leq$ " defining  $i_1 \leq i_2$  if and only if there exist  $a \in A_{i_1}, b \in A_{i_2}$ , such that  $b \circ a = b$ . This definition is consistent and the relation " $\leq$ " gives  $I$  the structure of a poset with the supremum property and if  $0 \in \tau(T)$  then there exists the greatest element  $1$  in  $I$ . For  $i \neq 1$  define  $\mathfrak{A}_i=(A_i; \{f_i^1\}_{i \in T^*})$  where  $f_i^1 = f_i|_{A_i}$  and  $\mathfrak{A}_1=(A_1; \{f_i^1\}_{i \in T})$  where  $f_i^1 = f_i|_{A_1}$ . If there are nullary fundamental operations in  $\mathfrak{A}$  then the values of all of them belong to  $A_1$ . It turns out that for  $i \neq 1$  any  $\mathfrak{A}_i$  is a subalgebra of  $\mathfrak{A}^*$  and  $\mathfrak{A}_1$  is a subalgebra of  $\mathfrak{A}$ . Finally define the mappings  $h_{i_1}^{i_2}: A_{i_1} \rightarrow A_{i_2}$  for  $i_1 \leq i_2$  by putting  $h_{i_1}^{i_2}(a) = a \circ b$  where  $b$  is an arbitrary element of  $A_{i_2}$ . The mappings so defined are  $T^*$ -homomorphisms and the system  $\mathcal{A} = \langle (I, \leq), \{\mathfrak{A}_i\}_{i \in I}, \{h_{i_1}^{i_2}\}_{i_1, i_2 \in I, i_1 \leq i_2} \rangle$  is a  $\tau$ -semilattice ordered system of algebras, for which  $\mathfrak{A}=\mathfrak{S}(\mathcal{A})$ . Conversely, every representation  $\mathfrak{A}=\mathfrak{S}(\mathcal{A})$  can be obtained by this construction starting with a suitable  $\tau$ -p-function  $\circ$  defined as follows; if  $a \in A_i, b \in A_j, u = \text{l.u.b.}(i, j)$  then  $a \circ b = h_i^u(a)$ . Moreover, the correspondence between  $\tau$ -p-functions of  $\mathfrak{A}$  and representations of  $\mathfrak{A}$  in the form  $\mathfrak{A}=\mathfrak{S}(\mathcal{A})$  is one-to-one.*

The proof of Theorem 2 is similar to that of Theorem 2 from [7]. The essential difference lies in this, that proving  $i \leq 1$  for any  $i \in I$  we use formula 4 for some  $t \in T_0$ . A  $\tau$ -p-function will be called an algebraic  $\tau$ -p-function if it is a realization of some binary term  $\varphi(x_1, x_2)$  not containing nullary fundamental operation symbols.

2. Let  $K$  be a variety of type  $\tau$ . We denote by  $E(K)$  the set of all identities satisfied in any algebra from  $K$ , by  $R_s(K)$  — the set of all strictly regular identities

from  $E(K)$ , by  $S_n(K)$  — the set of all nullary symmetrical identities from  $E(K)$ . Finally, let  $J(K) = R_s(K) \cup S_n(K)$ . We have

(viii) the set  $J(K)$  is closed under consequences (see [3]).

We denote by  $K_J$  the variety of type  $\tau$  defined by  $J(K)$ . We want to study the mapping  $K \rightarrow K_J$ . Let  $\mathbf{I}(\tau)$  denote the variety of type  $\tau$  defined by the set of all possible identities of type  $\tau$  being strictly regular or nullary symmetrical

(ix)  $\mathbf{I}(\tau)$  is not a trivial variety.

In fact consider an algebra  $\mathfrak{A} = (\{0, 1\}; \{f_t\}_{t \in T})$  where  $f_t = 1$  for any  $t \in T_0$  and if  $t \in T \setminus T_0$  then

$$f_t(n_1, \dots, n_{i(t)}) = \begin{cases} 1 & \text{if } 1 \in \{n_1, \dots, n_{i(t)}\} \\ 0 & \text{otherwise.} \end{cases}$$

Then  $\mathfrak{A} \in \mathbf{I}(\tau)$

(x) The variety  $\mathbf{I}(\tau)$  is equationally complete.

In fact any algebra  $\mathfrak{A} = (A; \{f_t\}_{t \in T}) \in \mathbf{I}(\tau)$  is roughly speaking a semilattice with the greatest element 1 if  $T_0 \neq \emptyset$ .

This follows from the fact that by assumption there exists a binary term, say  $x_1 \cdot x_2$  of type  $\tau$  and it satisfies in  $\mathbf{I}(\tau)$  the identities:  $x \cdot x = x$ ,  $x \cdot y = y \cdot x$ ,  $(x \cdot y) \cdot z = x \cdot (y \cdot z)$ ,  $f_{t_1} = f_{t_2}$  for any  $t_1, t_2 \in T_0$ ,  $x \cdot f_{t_1} = f_{t_1}$  for  $t_1 \in T_0$ ,  $f_t(x_1, \dots, x_{i(t)}) = x_1 \cdot \dots \cdot x_{i(t)}$  for any  $t \in T \setminus T_0$ . Moreover if  $\varphi = \psi$  is an identity such that no nullary fundamental operation symbol occurs in  $\varphi$  and some  $f_i$  occurs in  $\psi$  where  $t \in T_0$  then putting  $x$  for all variables in  $\varphi$  and  $\psi$  we get  $x = f_t$  and  $x = y$ . If no nullary fundamental operation symbols occur in  $\varphi$  and  $\psi$  and there exists  $x_i$  in  $\varphi$  which does not occur in  $\psi$  then putting  $x = x_i$  and  $y = x_j$  for  $j \neq i$  we get  $x = y$  or  $x \cdot y = y$ . Since  $x \cdot y = y \cdot x$  we get again  $x = y$ . Obviously, for any variety  $K$  we have

(xi)  $K_J = K \vee \mathbf{I}(\tau)$ .

The property (xi) together with the observation that  $\mathbf{I}(\tau)$  is practically the variety of semilattices explains us why we get the statement of Theorem 1. In fact a semilattice with 1 preserves only strictly regular and nullary symmetrical identities. Let  $U(K)$  denote the set of all equalities from  $E(K)$  in which nullary fundamental operation symbols do not occur and  $K_u$  be the variety of the type  $\tau^*$  defined by  $U(K)$ .

**COROLLARY 1.** *A variety  $K$  is defined only by strictly regular and nullary symmetrical identities, i.e.  $E(K) = J(K)$  iff  $\mathfrak{S}(\mathcal{A})$  belongs to  $K$  for any  $\iota$ -semilattice ordered system  $\mathcal{A}$  of algebras  $\mathfrak{A}_i$ , where  $\mathfrak{A}_i \in K_u$  for  $i \in I \setminus \{1\}$  and  $\mathfrak{A}_1 \in K$ .*

**PROOF.** The “only if” part follows from Theorem 1. The “if” part follows from the fact that the algebra  $\mathfrak{A}$  from (ix) belongs to  $\mathbf{I}(\tau)$ , generates  $\mathbf{I}(\tau)$  by (x) and  $\mathfrak{A}$  belongs to  $K$  since it is the sum of a  $\iota$ -semilattice ordered system of two 1-element algebras  $\mathfrak{A}_1 = (\{0\}, \{f_t\}_{t \in T^*})$  and  $\mathfrak{A}_2 = (\{1\}, \{f_t\}_{t \in T})$  where  $\mathfrak{A}_2 \in K$ ,  $\mathfrak{A}_1 \in K_u$ . In fact it is enough to put  $I = \{1, 2\}$   $1 \leq 2$  and to define the only  $T^*$ -homomorphism  $h$  different from identity putting  $h(0) = 1$ . Thus  $\mathbf{I}(\tau) \subseteq K$  and  $E(K) \subseteq J(\mathbf{I}(\tau))$ .

**THEOREM 3.** *Let  $K$  be a variety of algebras of the type  $\tau$  satisfying the following condition:*



(c) There exists a binary term  $x_1 \circ x_2$  of type  $\tau$  containing the variables  $x_1$  and  $x_2$ , containing no nullary fundamental operation symbols and such that the identity  $x_1 \circ x_2 = x_1$  belongs to  $E(K)$ .

Then an algebra  $\mathfrak{A}$  belongs to  $K_J$  iff  $\mathfrak{A} = \mathfrak{S}(\mathcal{A})$  for some 1-semilattice ordered system  $\mathcal{A}$  of algebras  $\mathfrak{A}_i$ , where  $\mathfrak{A}_i \in K_U$  for  $i \neq 1$  and  $\mathfrak{A}_1$  belongs to  $K$ .

PROOF.  $\leftarrow$  follows from Theorem 1.

$\rightarrow$  Since  $J(K) \subseteq E(\mathfrak{A})$  so the realization of  $x_1 \circ x_2$  in  $\mathfrak{A}$  is an algebraic 1- $p$ -function in  $\mathfrak{A}$ . Using Theorem 2 we conclude that  $\mathfrak{A} = \mathfrak{S}(\mathcal{A})$  for some 1-semilattice ordered system  $\mathcal{A}$  of algebras  $\mathfrak{A}_i$ . The proof that any  $\mathfrak{A}_i$  ( $i \in I$ ) satisfies all identities from  $U(K)$  is identical with that of Theorem 1 from [8]. Obviously,  $\mathfrak{A}_i$  satisfies all identities from  $S_n(K)$  as a subalgebra of  $\mathfrak{A}$ . It remains to prove that if an equality  $\varphi(x_{i_1}, \dots, x_{i_p}, c_{j_1}, \dots, c_{j_q}) = \psi(x_{k_1}, \dots, x_{k_r})$  belongs to  $E(K)$  where no nullary fundamental operation symbol occurs in  $\psi$  and some occur in  $\varphi$  then this equality is satisfied in  $\mathfrak{A}_i$ . But the equality  $\varphi = \psi \circ c_{j_i}$  belongs to  $J(K)$ . So for  $a_1, \dots, a_p, b_1, \dots, b_r \in A_i$  we have in  $\mathfrak{A}_i$

$$\varphi(a_1, \dots, a_p, c_{j_1}, \dots, c_{j_q}) = \psi(b_1, \dots, b_r) \circ c_{j_i}.$$

But  $\psi(b_1, \dots, b_r) \in A_i$  so  $\psi(b_1, \dots, b_r) \circ c_{j_i} = \psi(b_1, \dots, b_r)$ . Thus  $\varphi(a_1, \dots, a_p, c_{j_1}, \dots, c_{j_q}) = \psi(b_1, \dots, b_r)$  holds in  $\mathfrak{A}_i \in K$ . Q.E.D.

Let  $(L; +, \cdot)$  be a lattice. One says that for  $a, b \in L$   $a \neq b$  the element  $b$  covers  $a$  if for any  $c \in L$  we have  $a \leq c \leq b \Rightarrow c = a$  or  $c = b$ . Let  $L(\tau)$  be the lattice of all varieties of type  $\tau$ . It was proved in [1] that:

(xii) if  $A$  is an atom in  $L(\tau)$ ,  $V \in L(\tau)$ ,  $V$  does not contain  $A$  and for any algebra  $\mathcal{B} = (B, F) \in A \vee V$  we have  $\mathcal{B} \in V$  or there exists a congruence  $\theta \neq B \times B$  in  $\mathcal{B}$  such that  $\mathcal{B}/\theta \in A$  then  $A \vee V$  covers  $V$  in  $L(\tau)$ .

THEOREM 4. If  $K$  is a variety of type  $\tau$  satisfying the condition (c), then  $K_J$  covers  $K$  in the lattice  $L(\tau)$ .

PROOF. By (xi) and (x) we have  $K_J = I(\tau) \vee K$  and  $I(\tau)$  is an atom in  $L(\tau)$ . By assumption  $K \neq K_J$ . By Theorem 3 for any algebra  $\mathfrak{A} \in K_J$  we have  $\mathfrak{A} = \mathfrak{S}(\mathcal{A})$  where  $\mathcal{A}$  satisfies conditions of the statement of Theorem 3. Define in  $\mathfrak{A} = \mathfrak{S}(\mathcal{A})$  a relation  $\theta$  putting  $a \theta b$  if  $a$  and  $b$  belong to the same  $A_i$ . Then  $\theta$  satisfies assumption of (xii) and the theorem holds.

Let  $K$  be a variety of type  $\tau$  satisfying the condition (c). Assume that  $B$  is an equational base of  $K$  and  $U \subseteq B$  is an equational base of  $K_U$ . We define the set  $J'$  of identities as follows:

(a<sub>1</sub>)  $J'$  contains all strictly regular and all nullary symmetrical equalities from  $B$ .

(a<sub>2</sub>)  $J'$  contains identities (1)–(7) for  $t \in T$ .

(a<sub>3</sub>) if an equality  $\varphi(x_{i_1}, \dots, x_{i_p}) = \psi(x_{j_1}, \dots, x_{j_q})$  belongs to  $B$  where this equality is non regular and no nullary polynomial symbols occur in  $\varphi$  and  $\psi$ ,

$$\{x_{i_1}, \dots, x_{i_p}\} \setminus \{x_{j_1}, \dots, x_{j_q}\} = \{x_{k_1}, \dots, x_{k_r}\}$$

$$\{x_{j_1}, \dots, x_{j_q}\} \setminus \{x_{i_1}, \dots, x_{i_p}\} = \{x_{m_1}, \dots, x_{m_s}\}$$



then the equality

$$\varphi(x_{i_1}, \dots, x_{i_p}) \circ x_{m_1} \circ \dots \circ x_{m_s} = \psi(x_{j_1}, \dots, x_{j_q}) \circ x_{k_1} \circ \dots \circ x_{k_r}$$

belongs to  $J'$ ,

(a<sub>4</sub>) if  $\varphi(x_{i_1}, \dots, x_{i_p}, c_{j_1}, \dots, c_{j_q}) = \psi(x_{k_1}, \dots, x_{k_r})$  belongs to  $B$  where no nullary polynomial symbol occurs in  $\psi$  and some fundamental nullary operation symbols  $c_{j_1}, \dots, c_{j_q}$  occur in  $\varphi$  then the equality  $\varphi(x_{i_1}, \dots, x_{i_p}, c_{j_1}, \dots, c_{j_q}) = \psi(x_{k_1}, \dots, x_{k_r}) \circ c_{j_1}$  belongs to  $J'$ ,

(a<sub>5</sub>)  $J'$  is the smallest set satisfying (a<sub>1</sub>)—(a<sub>4</sub>).

**THEOREM 5.** *If  $K$  satisfies the above assumptions then  $J'$  is an equational base of  $K_J$ .*

**PROOF.** Denote by  $K'$  the variety of type  $\tau$  defined by  $J'$ . Obviously,  $K_J \subseteq K'$  since  $J' \subseteq J(K)$ . We shall prove that  $K' \subseteq K_J$ . By (a<sub>2</sub>)  $\circ$  is an algebraic  $p$ -function in any algebra  $\mathfrak{A} \in K'$ . By Theorem 2  $\mathfrak{A} = \mathfrak{S}(\mathcal{A})$  for some  $\iota$ -semilattice ordered system of algebras. But  $\mathfrak{A}_i$  is a subalgebra of  $\mathfrak{A}^*$  for  $i \neq \iota$ ,  $\mathfrak{A}_\iota$  is a subalgebra of  $\mathfrak{A}$  and any  $\mathfrak{A}_i$ ,  $i \in I$  satisfies  $x_1 \circ x_2 = x_1$  by Theorem 2. So any  $\mathfrak{A}_i$  satisfies all identities of  $U$  by (a<sub>1</sub>) and (a<sub>3</sub>) and  $\mathfrak{A}_\iota$  satisfies all identities of  $B$  by (a<sub>1</sub>), (a<sub>3</sub>) and (a<sub>4</sub>). Thus  $\mathfrak{A}_i \in K_U$  for  $i \neq \iota$ ,  $\mathfrak{A}_\iota \in K$ . Now, by Theorem 3,  $\mathfrak{A} \in K_J$  and consequently  $K' \subseteq K_J$ . Q.E.D

**COROLLARY 2.** *If  $K$  satisfies the assumptions of Theorem 5,  $K$  is finitely based,  $K_U$  is finitely based and  $|T|$  is finite then  $K_J$  is finitely based.*

Theorem 5 is a generalization of the result of E. Graczyńska, see [2].

The set  $J'$  so defined is in general too large and in practice we can choose a smaller set defining  $K_J$  as it shows the following example.

**EXAMPLE.** Let  $B$  be the variety of Boolean algebras with the fundamental operation symbols  $+$ ,  $\cdot$ ,  $'$ ,  $0$ ,  $1$ . Then the variety  $B_J$  can be described by the following system of identities:

- (1°)  $x + x = x \cdot x = x$
- (2°)  $x + y = y + x, \quad x \cdot y = y \cdot x$
- (3°)  $(x + y) + z = x + (y + z), \quad (x \cdot y) \cdot z = x \cdot (y \cdot z)$
- (4°)  $(x + y) \cdot z = x \cdot z + y \cdot z, \quad (x \cdot y) + z = (x + z) \cdot (y + z)$
- (5°)  $(x')' = x$
- (6°)  $(x + y)' = x' \cdot y'$
- (7°)  $x + x \cdot x' = x$
- (8°)  $x \cdot x' + y \cdot y' = x \cdot x' \cdot y \cdot y'$
- (9°)  $x \cdot 0 = 0$
- (10°)  $x \cdot x' + 0 = 0$
- (11°)  $0' = 1.$

In fact, let us denote by  $D$  the variety of the same type as  $B$  defined by  $(1^\circ)$ — $(11^\circ)$ . Obviously,  $B_J \subseteq D$  since all axioms  $(1^\circ)$ — $(11^\circ)$  belong to  $J(B) = E(B_J)$ . Let  $\mathfrak{A} = (A; +, \cdot, ', 0, 1) \in D$ . We show that the realization of the term  $x_1 \circ x_2 = x_1 + x_1 \cdot x_2$  in  $\mathfrak{A}$  is a  $\iota$ - $p$ -function. It was shown in [9] that  $\circ$  satisfies (1)—(3) and (4)—(7) for the operations  $+$ ,  $\cdot$  and  $'$ . By  $(9^\circ)$  we have  $0 \circ x = 0 + 0 \cdot x = 0$  so (4) holds for 0. By  $(6^\circ)$  and by  $(11^\circ)$  we have  $1 \circ x = 1 + 1 \cdot x = (0 + 0 \cdot x)' = 0' = 1$ . So (4) holds for 1. Thus  $\circ$  is a  $\iota$ - $p$ -function in  $\mathfrak{A}$ . By Theorem 2,  $\mathfrak{A} = \mathfrak{S}(\mathcal{A})$  for some  $\iota$ -semi-lattice ordered system  $\mathcal{A}$  of algebras  $\mathfrak{A}_i$ . By results of [9] any algebra  $(A_i; +, \cdot, ', 0)$  ( $i \in I$ ) is a Boolean algebra. We show that for  $a \in A_i$ ,  $0 = a \cdot a'$  and  $1 = a + a'$ . Since  $a \cdot a'$ ,  $0 \in A_i$  we have by Theorem 2 and  $(10^\circ)$ :

$$(a \cdot a') \circ 0 = a \cdot a' \quad \text{and} \quad (a \cdot a') \circ 0 = (a \cdot a') + a \cdot a' \cdot 0 = a \cdot a' + 0 = 0.$$

Also  $1 = 0' = (a \cdot a')' = a + a'$  by  $(11^\circ)$  and  $(6^\circ)$ .

Since Boolean algebras are equationally complete so  $\mathcal{A}$  is such that  $\mathfrak{A}_i = (A_i; +, \cdot, ', 0) \in B_U$  for  $i \neq \iota$  and  $\mathfrak{A}_\iota = (A_\iota; +, \cdot, ', 0, 1) \in B$ . Now by Theorem 3:  $\mathfrak{A} = \mathfrak{S}(\mathcal{A}) \in K_J$  and  $D \subseteq K_J$  what ends the proof.

REMARK. Many important varieties satisfy assumptions of Theorem 3 and 4. We have  $x_1 \cdot x_2 \cdot x_2^{-1} = x_1$  in groups,  $x_1 + x_1 \cdot x_2 = x_1$  in lattices, Boolean algebras. So we can apply Theorem 3 and 4 for this varieties.

# REFERENCES

- [1] DUDEK, J. and PŁONKA, J., On covering in lattices of varieties of algebras, *Bull. Acad. Polon. Sci. Ser. Math.* **31** (1983), 1—4.
- [2] GRACZYŃSKA, E., On regular identities, *Algebra Universalis* **17** (1983), 369—375.
- [3] GRÄTZER, G., *Universal Algebra*, Springer-Verlag, Berlin, 1979. MR 80g: 08001.
- [4] JOHN, R., On classes of algebras definable by regular equations, *Colloq. Math.* **36** (1976), 17—21. MR 54# 12604.
- [5] LAKSER, H., PADMANABHAN, R. and PLATT, C. R., Subdirect decomposition of Płonka sums, *Duke Math. J.* **39** (1972), 485—488. MR 46# 3416.
- [6] MITSCHKE, A., On a representation of groupoids as sums of direct systems, *Colloq. Math.* **28** (1973), 11—18. MR 48# 10953.
- [7] PŁONKA, J., On a method of construction of abstract algebras, *Fund. Math.* **61** (1967), 183—189. MR 37# 1294.
- [8] PŁONKA, J., On equational classes of abstract algebras defined by regular equations, *Fund. Math.* **64** (1969), 241—247. MR 39# 5450.
- [9] PŁONKA, J., On sum of direct systems of Boolean algebras, *Colloq. Math.* **20** (1969), 209—214. MR 40# 2590.
- [10] PŁONKA, J., On the sum of a direct system of universal algebras with nullary polynomials, *Algebra Universalis* **19** (1984), 197—207.

( Received September 8. 1983 )

UL. ŚLEŻNA 9, M. 6  
PL—53-301 WROCŁAW  
POLAND



## SOME RESULTS ON THE CANTOR SET

HARRY I. MILLER<sup>1</sup> and POLYCHRONIS J. XENIKAKIS

### Abstract

It is shown that  $C$ , the Cantor set, contains a subset similar (in the sense of elementary geometry) to each three point subset of the real line, that  $C$  contains no arithmetic progression of length five, and there exists a four point subset of the reals that is not similar to any subset of  $C$ . A generalization of the last mentioned result is also presented.

### 1. Introduction

All sets considered in this paper are subsets of  $\mathbb{R}$  the real line. If  $A$  is Lebesgue measurable and  $mA$ , the Lebesgue measure of  $A$ , is positive then by a classical result of Steinhaus [5],  $D(A)$  contains a non-empty interval, where  $D(A)$  is defined by the equation  $D(A) = \{x - y : x, y \in A\}$ . The set  $A$  is called a Baire set if  $A$  can be written in the form  $A = (G \setminus P) \cup Q$ , where  $G$  is an open set and  $P$  and  $Q$  are sets of the first Baire category. Piccard [3] proved a Baire set analogue of the result of Steinhaus; namely she proved that  $D(A)$  contains a non-empty interval if  $A$  is a Baire set of the second Baire category.

Sets  $M$  and  $N$  will be called similar in case there exists a function  $f, f: \mathbb{R} \rightarrow \mathbb{R}$ , of the form  $f(x) = ax + b$ ,  $a \neq 0$ , such that  $f(M) = N$ . Ruziewicz [4] showed that if  $mA > 0$  and  $E$  is any finite set, then  $A$  contains a subset  $E'$ , such that  $E$  and  $E'$  are similar. Miller [2] proved the Baire set analogue of this result.

It is natural to consider the following classes.

$$\mathcal{A} = \{A : D(A) \text{ contains a non-empty interval}\}.$$

$$\mathcal{B} = \{B : \text{for each finite set } E, B \text{ contains a subset } E', \text{ such that } E \text{ and } E' \text{ are similar}\}.$$

If we set

$$\mathcal{M}^+ = \{A : A \text{ is measurable and } mA > 0\} \text{ and}$$

$$\mathcal{Ba}^+ = \{A : A \text{ is a Baire set of the second Baire category}\}$$

then by the four previously mentioned theorems we have :

$$\mathcal{M}^+ \cup \mathcal{Ba}^+ \subseteq \mathcal{A} \cap \mathcal{B}.$$

<sup>1</sup> The work of the first author was supported by the Scientific Research Fund of Bosna and Herzegovina.

1980 *Mathematics Subject Classification*. Primary 04A15; Secondary 28A05.

*Key words and phrases*. Cantor set and geometric similarities, Cantor set and arithmetic progressions.

In a paper recently submitted for publication the present authors have proved several results about these four classes.

Let  $C$  denote the Cantor set. Since  $C$  is a first category set of Lebesgue measure zero  $C \notin \mathcal{M}^+ \cup \mathcal{B}a^+$ , even though it is Lebesgue measurable and a Baire set. However, it is well-known that  $D(C) = [-1, 1]$  and therefore  $C \in \mathcal{A}$ . An even stronger result about  $C$  is true. Boas [1] has shown that for almost all  $d$  in  $[-1, 1]$ , the set  $\{(x, y) : x, y \in C, x - y = d\}$  has the cardinal number of the continuum. What about  $\mathcal{B}$ ? Is  $C \in \mathcal{B}$  or is  $C \notin \mathcal{B}$ ? To the best of our knowledge this question has not been answered in the literature. In this article we settle this question, showing that  $C \notin \mathcal{B}$ . Professor F. Galvin, on reading a handwritten copy of our results on the Cantor set, has suggested a generalization of our results which has been adapted in Theorem 1 of this paper.

## 2. Results

The Cantor set is the disjoint union of two sets each of which is similar to the entire Cantor set. We now present a result which shows, as a corollary, that  $C \notin \mathcal{B}$ .

**THEOREM 1.** *Suppose that  $A = A_1 \cup A_2$ ;  $A$ ,  $A_1$ , and  $A_2$  are similar and that  $\sup A_1 < \inf A_2$ . Then  $A \notin \mathcal{B}$ .*

**PROOF.** Since  $A$ ,  $A_1$ , and  $A_2$  are similar and  $A_1$  is bounded above and  $A_2$  is bounded below it follows at once that  $A$  is bounded. For convenience suppose that  $\inf A = 0$  and  $\sup A = 1$ . Let  $d = \inf A_2 - \sup A_1$ . Then  $d > 0$ . We will now show, using the "nested interval theorem", that  $A$  contains no subset similar to  $\{0, 1, 2, \dots, n\}$  if  $n$  is sufficiently large. Suppose that  $\{x_0, x_1, x_2, \dots, x_n\} \subseteq A$ ;  $x_i < x_{i+1}$  for each  $i$  and that  $\{x_0, x_1, x_2, \dots, x_n\}$  is similar to  $\{0, 1, \dots, n\}$ . We will now show that all of these points either lie in  $A_1$  or  $A_2$  if  $n$  is sufficiently large. This is clear, for if the points  $x_0, x_1, x_2, \dots, x_n$  are not all in  $A_1$  or  $A_2$  then there exists an  $i_0$  such that  $x_{i_0} \in A_1$  and  $x_{i_0+1} \in A_2$ . Therefore  $x_{i_0+1} - x_{i_0} \geq d$  and hence

$$x_n - x_0 = \sum_{i=0}^{n-1} (x_{i+1} - x_i) \geq nd \quad \text{since } \{0, 1, \dots, n\}$$

and  $\{x_0, x_1, \dots, x_n\}$  are similar. If  $n$  is chosen so that  $nd > 1$ , then we get  $x_n - x_0 > 1$  which is impossible. Let  $A^1$  denote the set (either  $A_1$  or  $A_2$ ) that contains all the points  $x_0, x_1, x_2, \dots, x_n$ . Since  $A^1$  is similar to  $A$ , we have  $A^1 = A_1^1 \cup A_2^1$  and  $\inf A_2^1 - \sup A_1^1 = |a|d$ , where  $f(x) = ax + b$  is the similarity map between  $A$  and  $A^1$ . We now claim that the points  $x_0, x_1, x_2, \dots, x_n$  all lie in  $A_1^1$  or  $A_2^1$ . If this is not the case then there exists an  $i_0$  such that  $x_{i_0} \in A_1^1$  and  $x_{i_0+1} \in A_2^1$  and therefore  $x_{i_0+1} - x_{i_0} \geq |a|d$  which implies

$$x_n - x_0 \geq n|a|d > |a| = \sup A^1 - \inf A^1$$

which is impossible since  $\{x_0, x_1, x_2, \dots, x_n\} \subseteq A^1$ .

Continuing, we get  $\{x_0, x_1, x_2, \dots, x_n\} \subseteq A^m$  for each natural number  $m$ , where  $A^m$  is one of the two similar sets whose union is  $A^{m-1}$ . Since  $\{A^m\}$  is a nested sequence of sets whose diameters tend to zero it follows that  $A$  contains no subset similar to  $\{0, 1, 2, \dots, n\}$  if  $nd > 1$ . Therefore  $A \notin \mathcal{B}$ .

For  $C$ , the Cantor set, we have  $C = \{C \cap [0, 1/3]\} \cup \{C \cap [2/3, 1]\}$  and  $d = 1/3$ . Therefore by the argument in the proof of Theorem 1 the following is true.

**COROLLARY 2.**  $C$ , the Cantor set, contains no arithmetic progression of length five and therefore  $C \notin \mathcal{B}$ .

**REMARK 1.**  $C$  does contain an arithmetic progression of length four, namely  $0, 1/3, 2/3, 1$ . However, a slight adjustment of the argument used in the proof of Theorem 1 shows that  $C$  contains no subset similar to the set  $\{0, 1, 2, 3^{1/4}\}$ . Namely, if  $\{x_0, x_1, x_2, x_3\}$  is similar to  $\{0, 1, 2, 3^{1/4}\}$  and  $\{x_0, x_1, x_2, x_3\} \subseteq C$ , then it follows easily that  $\{x_0, x_1, x_2, x_3\}$  is either a subset of  $[0, 1/3]$  or of  $[2/3, 1]$ . Continuing one shows that  $\{x_0, x_1, x_2, x_3\}$  must be contained in an interval of length  $3^{-m}$  for each  $m$  which is impossible.

We now show that  $C$  contains a set similar to each three point set.

**THEOREM 3.**  $C$  contains a subset similar to each three point set.

**PROOF.** Utz [6] showed that for any line:  $y - mx = c$ , with  $1/3 \leq |m| \leq 3$  and  $0 \leq c \leq 1$ , there exist numbers  $x_1, y_1$  both in the Cantor set such that  $(x_1, y_1)$  lies on the line, i.e.  $y_1 - mx_1 = c$ .

Suppose  $S = \{s_1, s_2, s_3\}$  and  $s_1 < s_2 < s_3$ . Let  $q = (s_2 - s_1)/(s_3 - s_2)$ . Then  $q > 0$ . We will now find a subset  $C'$  of  $C$  that is similar to  $S$ .

*Case I.* Suppose  $1/3 \leq q \leq 3$ . We will prove that for each such  $q$  there exists  $x, c, y \in C$ , with  $x \neq y$ , such that  $y = (c + c/q) - x/q$  (a). To see this consider the following sketch.

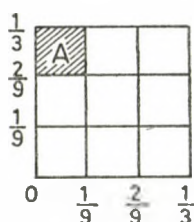


Fig. 1

For  $1/3 \leq q \leq 1$ , pick  $c = 1/9$ . A simple check shows that the line (a) intersects the square  $A$  and therefore by the argument in Utz's paper there exists  $x, y \in C$  with  $(x, y)$  in the square  $A$  and such that  $(x, y)$  lies on line (a). Therefore  $x \neq y$ . For  $1 < q \leq 3$  pick  $c = 2/9$ . Again a simple check shows that the line (a) intersects square  $A$  and therefore by Utz's argument there exists  $(x, y) \in C \times C$ , with  $(x, y)$  in square  $A$  such that  $(x, y)$  lies on line (a). Therefore  $x \neq y$ .

In either case  $x \neq c$  (as  $x = c$  implies  $y = c$  which implies  $x = y$ ) and also  $y \neq c$  (as  $y = c$  implies  $x = c$  which implies  $x = y$ ). Therefore  $(c - x)/(y - c) = q$ . From this it follows that  $\{x, c, y\}$  and  $\{s_1, s_2, s_3\}$  are similar.

*Case II.* Suppose  $q > 3$ . Consider (a), with  $x = 0$ . Then we have  $y = c(1/q + 1)$ . Since  $q > 3$ ,  $1/q + 1 \in (1, 4/3)$ . Again using Utz's theorem, with  $c = 0$ , we see that each number  $m$ , satisfying  $1/3 \leq |m| \leq 3$ , can be written in the form  $m = y/x$ , with  $x, y \in C$ . Since  $(1, 4/3) \subseteq [1/3, 3]$ , each number  $(1/q + 1)$ , with  $q > 3$ , can be written



as  $y/c$ , with  $y, c \in C$ . Therefore  $0, c, y \in C$ ,  $0 < c < y$  and  $(c-0)/(y-c) = q$ , or  $\{0, c, y\}$  and  $\{s_1, s_2, s_3\}$  are similar.

*Case III.* Suppose  $0 < q < 1/3$ . This case reduces to case II, for if  $(s_2 - s_1)/(s_3 - s_2) < 1/3$ , then  $(s'_2 - s'_1)/(s'_3 - s'_2) > 3$ , where  $s'_1 = -s_3$ ,  $s'_2 = -s_2$ , and  $s'_3 = -s_1$ .

**REMARK 2.** We now present an alternative proof of the fact that  $C$  contains no five point set in arithmetic progression, that is for each  $a$  and  $d$ ,  $d \neq 0$ ,  $\{a, a+d, a+2d, a+3d, a+4d\} \not\subseteq C$ .

**PROOF.** First we will show that if  $x, y, (x+y)/2$  are all in  $C$  and  $x \neq y$ , then there exists a unique natural number  $n$  such that  $|x-y| = 2/3^n$ . To see this suppose  $x, y, (x+y)/2$  are all in  $C$ ,  $x \neq y$  and that  $x = 0.x_1x_2x_3 \dots$  and  $y = 0.y_1y_2y_3 \dots$  are ternary developments of  $x$  and  $y$  that only use the digits 0 and 2. Since  $x \neq y$  there exists a natural number  $i$  such that  $x_i \neq y_i$ . Let  $n$  be the smallest such number. Then  $x_n + y_n = 2$ . Furthermore  $(x+y)/2 = 0.z_1z_2z_3 \dots$ , where  $z_m = (x_m + y_m)/2$  for each  $m$ . Therefore  $z_n = 1$ . If there exists  $n' > n$  such that  $x_{n'} \neq y_{n'}$ , then  $z_{n'} = 1$  which implies that the development of  $(x+y)/2$  has a second 1. This is a contradiction of the fact that  $(x+y)/2 \in C$ . Therefore if  $x, y, (x+y)/2$  all lie in  $C$  and  $x \neq y$ , then there exists a unique  $n$  such that  $x_n \neq y_n$ . Then clearly  $|x-y| = 2/3^n$ . Suppose that  $C$  contains five successive terms of an arithmetic progression, i.e.  $\{a, a+d, a+2d, a+3d, a+4d\} \subseteq C$  for some  $a$  and  $d$ ,  $d \neq 0$ . Then, by the above argument we have  $|(a+4d)-a| = 2/3^n$  for some natural number  $n$  and  $|(a+3d)-(a+d)| = 2/3^k$  for some natural number  $k$ . This implies that  $2/3^n = 2 \cdot 2/3^k$  or that  $3^{k-n} = 2$ , which is impossible. Therefore  $C$  cannot contain five successive terms of an arithmetic progression.

**REMARK 3.** If  $A = A_1 \cup A_2$ ;  $A$  is closed;  $A, A_1$  and  $A_2$  are similar and  $\sup A_1 < \inf A_2$ , then it is a straightforward exercise to show that  $mA = 0$  and that  $A$  is a set of the first Baire category. These facts also follow from Theorem 1 and the fact (mentioned in the introduction) that  $\mathcal{M}^+ \cup \mathcal{B}a^+ \subseteq \mathcal{A} \cap \mathcal{B}$ .

#### REFERENCES

- [1] BOAS, R. P., The distance set of the Cantor set, *Bull. Calcutta Math. Soc.* **54** (1962), 103—104.
- [2] MILLER, H. I., Relationships between various gauges of the size of sets of real numbers (II), *Akad. Nauka; Umjet. Bosne; Hercegov. Rad. Odjelj. Prirod. Mat. Nauka*, **59**, Knjiga 16 (1976), 37—48.
- [3] PICCARD, S., *Sur les ensembles de distances des points d'un espace Euclidien*, Neuchâtel, 1933.
- [4] RUZIEWICZ, M. S., Contribution à l'étude des ensembles de distances de points, *Fund. Math.* **7** (1925), 141—143.
- [5] STEINHAUS, H., Sur les distances des points des ensembles de mesure positive, *Fund. Math.* **1** (1920), 93—104.
- [6] UTZ, W. R., The distance set for the Cantor discontinuum, *Amer. Math. Monthly* **58** (1951), 407—408.

(Received September 15, 1983)

ODSJEK ZA MATEMATIKU  
PRIRODNO-MATEMATICKOG FAKULTETA  
UNIVERZITETA U SARAJEVU  
VOJVODE PUTNIKA 43  
YU—71000 SARAJEVO  
YUGOSLAVIA  
DEPARTMENT OF MATHEMATICS  
ARISTOTLE UNIVERSITY OF THESSALONIKI  
THESSALONIKI  
GREECE



# О ФУНДАМЕНТАЛЬНОЙ ПРИВОДИМОСТИ САМОСОПРЯЖЕННЫХ И УНИТАРНЫХ ОПЕРАТОРОВ В ПРОСТРАНСТВАХ С ИНДЕФИНИТНОЙ МЕТРИКОЙ

Ц. БАЯСГАЛАН

Настоящая работа является непосредственным продолжением работы [1]. В работе исследуется фундаментальная приводимость самосопряженных и унитарных операторов. Основные результаты содержатся в диссертации автора [2].

Фундаментальная приводимость операторов разных классов изучалась в работах Р. Филлипса [3], Е. Песонена [4], Р. Кюне [5], П. Хесса [6] и Б. Мэк Энниса [13].

Автор выражает искреннюю благодарность П. Йонасу [Берлин] за ценный совет. В частности, он обратил мое внимание к работе Р. В. Акопяна [7].

## §1. Некоторые определения и предложения из теории пространств с индефинитной метрикой

Необходимые определения и предложения изложены в книге Я. Богнара [8].

**Теорема 1.1** ([8], стр. 108). *Подпространство в пространстве Крейна ортодополняемо тогда и только тогда, когда оно есть ортогональная прямая суммы замкнутого равномерно положительного и замкнутого равномерно отрицательного подпространств.*

**Теорема 1.2** ([8], стр. 43). *Пусть  $(\cdot, \cdot)$  — индефинитное внутреннее произведение на линейном пространстве  $H$ ,  $a(\cdot, \cdot)_1$  — другое внутреннее произведение на том же пространстве. Если из  $(x, x) = 0$  следует  $(x, x)_1 = 0$ , то для некоторого вещественного числа  $\mu_1$  имеем  $(x, y)_1 = \mu_1(x, y)$  для всех  $x, y$  из  $H$ .*

В дальнейшем рассмотрим только ограниченные линейные операторы, определенные на индефинитном пространстве Крейна  $H$ . Пусть  $A$  — самосопряженный оператор в  $H$ . Оператор  $A$  называется равномерно положительным, если имеет место соотношение

$$(Ax, x) \geq \alpha \|x\|_J^2 \quad (x \in H),$$

где  $\alpha$  — положительное постоянное,  $\|\cdot\|_J$  — некоторая  $J$ -норма. Оператор  $A$  называется строго положительным, если  $(Ax, x) > 0$  при  $x \neq 0$ ,  $x \in H$ .

1980 *Mathematics Subject Classification*. Primary 47B50.

*Key words and phrases*. Krein space, selfadjoint operators.

Теорема 1.3 ([8], стр. 122). Для произвольного оператора  $A$  имеет место соотношение  $A^+ = JA^*J$ , где  $A^+$ —сопряженный оператор к  $A$  относительно  $J$ -скалярного произведения, а  $A^*$ —сопряженный оператор к  $A$  относительно индефинитной метрики.

Теорема 1.4 ([8], стр. 136). Пусть  $H$ —пространство Крейна.

а) Пусть  $A$ —самосопряженный оператор и не вещественное число  $\alpha$  не принадлежит спектру оператора  $A$ . Тогда для любого  $\mu$  с  $|\mu|=1$  оператор

$$(1.1) \quad U = \mu(A - \bar{\alpha}I)(A - \alpha I)^{-1} = \mu(I + (\alpha - \bar{\alpha})(A - \alpha I)^{-1})$$

является унитарным и  $\mu$  не принадлежит спектру оператора  $U$ .

б) Пусть  $U$ —унитарный оператор,  $|\mu|=1$ , и  $\mu$  не принадлежит спектру оператора  $U$ . Тогда для любого не вещественного числа  $\alpha$  оператор

$$(1.2) \quad A = (\alpha U - \bar{\alpha}\mu I)(U - \mu I)^{-1} = \alpha I + \mu(\alpha - \bar{\alpha})(U - \mu I)^{-1}$$

является самосопряженным и  $\alpha$  не принадлежит спектру оператора  $A$ .

Преобразования (1.1) и (1.2) взаимно обратны; они называются преобразованиями Кэли.

Оператор  $T$  называется фундаментально приводимым, если существует такое фундаментальное разложение  $H = H^+ (+) H^-$ , что  $TH^+ \subset H^+$ ,  $TH^- \subset H^-$ .

Теорема 1.5 ([3]). Унитарный оператор  $U$  в пространстве  $H$  фундаментально приводим в том и только в том случае, если последовательность  $\{\|U^n\|_J\}_{n=1}^\infty$  ограничена.

Теорема 1.6 ([3]). Для коммутативной группы  $\mathfrak{U}$  унитарных операторов в пространстве  $H$ , ограниченной по  $J$ -норме, существует фундаментальное разложение, приводящее одновременно все операторы из этой группы.

## §2. Фундаментальная приводимость самосопряженных и унитарных операторов

Сначала мы сделаем некоторые дополнения к работе [1].

Лемма 2.1. Пусть  $A$  положительный оператор в  $H$ . Тогда спектр  $A$  содержит неположительные, так и неотрицательные значения.

Доказательство. Рассмотрим случай, когда  $0 \notin \sigma(A)$ , где  $\sigma(A)$ —спектр оператора  $A$ . Тогда по теореме 1 из [1] оператор  $A$  фундаментально приводим. Следовательно, сужения оператора  $A$  на соответствующие компоненты фундаментального разложения имеют положительный и отрицательный спектр.

Тем самым не существует положительный оператор с положительным спектром в индефинитном пространстве Крейна.

Лемма 2.2. Оператор  $A$  равномерно положителен в том и только в том случае, если он является положительным оператором, имеющим ограниченный обратный.

Доказательство. Необходимость вытекает из леммы VII. 3.2 из [8].

Обратно, по теореме 1 из [1] оператор  $A$  фундаментально приводим. Следовательно, для некоторой фундаментальной симметрии  $J$  имеем

$$\inf_{\|x\|_J=1} (Ax, x) = \inf_{\|x\|_J=1} (JAx, x)_J > 0$$

в силу того, что оператор  $JA$  является  $J$ -самосопряженным, имеющим ограниченный обратный.

Лемма 2.3. Пусть  $\{T_\alpha\}$ —семейство операторов, коммутирующих с положительным оператором  $A$ , имеющим ограниченный обратный. Тогда существует фундаментальное разложение, приводящее одновременно все операторы  $T_\alpha$ .

Доказательство. Оператор  $A$  фундаментально приводим. Рассмотрим разложение (см. доказательство теоремы 3 из [1])

$$H = E(-0)H(\dot{+})(I - E(+0))H(\dot{+})(E(+0) - E(-0))H,$$

где подпространство  $(E(+0) - E(-0))H$  есть ядро оператора  $A$ . Тогда легко видеть, что фундаментальное разложение

$$H = E(-0)H(\dot{+})(I - E(+0))H$$

приводит каждый оператор  $T_\alpha$ .

Критерий регулярности спектральной функции (т.е. существования сильных пределов  $E(+0)$  и  $E(-0)$ ) приведен в работе Р. В. Акопяна [7]. Полученные в этой работе результаты дают критерий фундаментальной приводимости положительного оператора в терминах его резольвенты.

Теорема А ([7]). Пусть  $A$ —строго положительный оператор в  $H$ . Тогда спектральная функция оператора  $A$  регулярна тогда и только тогда, когда

$$\|(A - iy)^{-1}\|_J = o\left(\frac{1}{|y|}\right) \quad \text{при } y \rightarrow 0,$$

где  $\|\cdot\|_J$ —некоторая  $J$ -норма.

Теорема Б ([7]). Пусть  $A$ —положительный оператор в  $H$ . Спектральная функция оператора  $A$  регулярна тогда и только тогда, когда существует такое постоянное  $C$ , что неравенство

$$\|(A - iy)^{-k}\|_J \leq \frac{C}{y^k} \quad (y > 0, k = 1, 2, \dots)$$

выполняется на ядре оператора  $S$  (где  $S$ —оператор, участвующий в спектральном представлении

$$A = S + \int_{-\infty}^{\infty} t dE(t).$$

Следующие две теоремы непосредственно вытекают из теорем А и Б, если привлечь теорему 3 из [1].

**Теорема 2.4.** Пусть  $A$  — строго положительный оператор в  $H$ . Оператор  $A$  фундаментально приводим тогда и только тогда, когда

$$\|(A - iy)^{-1}\|_J = o\left(\frac{1}{|y|}\right) \quad \text{при } y \rightarrow 0.$$

**Теорема 2.5.** Пусть  $A$  — положительный оператор в  $H$ . Оператор  $A$  фундаментально приводим в том и только том случае, если

- а) ядро оператора  $A$  ортогодополняемо,
- б)

$$\|(A - iy)^{-k}\|_J \leq \frac{C}{y^k} \quad (y > 0, k = 1, 2, \dots),$$

где  $C$  — положительное постоянное.

Теперь мы приходим к рассмотрению фундаментальной приводимости самосопряженных и унитарных операторов в пространстве Крейна.

**Теорема 2.6.** а) Самосопряженный оператор  $A$  фундаментально приводим в том и только том случае, если для некоторой фундаментальной симметрии  $J$ , оператор  $A$  подобен  $J$ -самосопряженному оператору;

б) самосопряженный оператор  $A$  фундаментально приводим в том и только том случае, если для некоторой фундаментальной симметрии  $J$  он является  $J$ -самосопряженным оператором.

Доказательство этой теоремы аналогично доказательству теоремы VIII. 1.4 из [8].

Из этой теоремы вытекает, что спектр фундаментально приводимого самосопряженного оператора вещественен и его собственные значения полупросты. Заметим также, что ядро фундаментально приводимого оператора ортогодополняемо.

**Теорема 2.7.** Для самосопряженного оператора  $A$  следующие два условия эквивалентны:

- а)  $A$  — фундаментально приводимый оператор;
- б) спектр оператора  $A$  вещественен и

$$\|(A - z)^{-1}\|_J \leq \frac{1}{|\operatorname{Im} z|} \quad (\operatorname{Im} z \neq 0),$$

где  $\|\cdot\|_J$  — некоторая  $J$ -норма.

**Доказательство.** Если оператор  $A$  фундаментально приводим, то по предыдущей теореме он является самосопряженным оператором в некотором  $J$ -скалярном произведении, следовательно, по теории гильбертовых пространств верен пункт б).

Обратно, при выполнении условия б) оператор  $A$  является самосопряженным оператором в  $J$ -скалярном произведении, в силу чего, он фундаментально приводим.

Используя одну идею работы [7] докажем следующую теорему.

**Теорема 2.8.** Для самосопряженного оператора  $A$  следующие два условия эквивалентны:

- а)  $A$  — фундаментально приводим;
- б) спектр оператора  $A$  вещественен и

$$\|(A - i\gamma)^{-k}\|_J \leq \frac{C}{|\gamma|^k} \quad (k = 1, 2, \dots, \gamma \neq 0, \gamma \in \mathbb{R}),$$

где  $\|\cdot\|_J$  — некоторая  $J$ -норма,  $C$  — положительное постоянное.

**Доказательство.** Пусть оператор  $A$  фундаментально приводим. Тогда пункт б) вытекает из теоремы 2.7.

Пусть условие б) выполнено. Тогда

$$\|(i\gamma^{-1}A + 1)^{-k}\|_J \leq C \quad (k = 1, 2, \dots, \gamma \neq 0, \gamma \in \mathbb{R}).$$

Поэтому по теории полугрупп операторов (см. [9], стр. 347) существует такое постоянное  $C_1$ , что

$$\|e^{iAt}\|_J \leq C_1 \quad (t \in \mathbb{R}).$$

По теореме Секефальви-Надя (см. [10]), существует скалярное произведение  $(\cdot)_1$ , эквивалентное  $J$ -скалярному произведению, что операторы  $e^{iAt}$  являются унитарными в  $(\cdot)_1$ . Легко видно, что оператор  $A$  является самосопряженным в  $(\cdot)_1$ . Поэтому преобразование Кэли  $U$  оператора  $A$  есть унитарный оператор в  $(\cdot)_1$ . Тем самым  $\|U^n\|_1 = 1$  ( $n = 1, 2, \dots$ ), где  $\|\cdot\|_1$  — соответствующая норма. В исходной  $J$ -норме имеет место неравенство  $\|U^n\|_J \leq C_2$  ( $n = 1, 2, \dots$ ), где  $C_2$  — некоторое постоянное.

По теореме 1.5 оператор  $U$  фундаментально приводится некоторым фундаментальным разложением  $H = H^+(\frac{1}{2})H^-$ . Его преобразование Кэли (см. теорему 1.4) оператор  $A$  приводится этим же разложением.

Теперь рассмотрим один критерий фундаментальной приводимости для изометрического и унитарного оператора.

**Теорема 2.9.** Пусть  $U$  — изометрический оператор в  $H$ . Оператор  $U$  фундаментально приводим тогда и только тогда, когда

$$\|U\|_J \leq 1,$$

где  $\|\cdot\|_J$  — некоторая  $J$ -норма.

**Доказательство.** Пусть  $U$  фундаментально приводим. Тогда мы имеем  $UJ = JU$ , где  $J$  — некоторая фундаментальная симметрия. Далее,

$$(Ux, Uy)_J = (JUx, Uy) = (UJx, Uy) = (Jx, y) = (x, y)_J.$$

Тем самым, оператор  $U$  является изометрией в  $J$ -скалярном произведении. Поэтому  $\|U\|_J \leq 1$ .

Обратно, пусть  $\|U\|_J \leq 1$ . Рассмотрим соответствующее фундаментальное разложение  $H = H^+(\frac{1}{2})H^-$ . Тогда либо  $UH^+ \subset H^+$ , либо существует вектор  $x^+ \in H^+$  с  $Ux^+ \notin H^+$ .

Докажем, что последний случай невозможен. Мы можем предполагать,

что  $\|x^+\|_J = 1$ . Положим  $Ux^+ = y^+ + y^-$ , где  $y^+ \in H^+$ ,  $y^- \in H^-$ . По условию  $y^- \neq 0$ . Следовательно,

$$\|Ux^+\|_J^2 = \|y^+\|_J^2 + \|y^-\|_J^2 > \|y^+\|_J^2 - \|y^-\|_J^2 = (Ux^+, Ux^+) = (x^+, x^+) = \|x^+\|_J^2 = 1.$$

Отсюда  $\|U\|_J > 1$ , а это противоречие. Значит  $UH^+ \subset H^+$ .

Аналогично доказывается включение  $UH^- \subset H^-$ .

Теперь сформулируем некоторые критерии фундаментальной приводимости самосопряженных и унитарных операторов в пространстве Понтрягина.

**Лемма 2.10.** Пусть  $A$  — самосопряженный оператор в  $H$ . Если  $\text{Ker } A$  ортодополняемо и  $M$  — ортодополняемое подпространство, инвариантное относительно оператора  $A$ , то  $\text{Ker } (A|M)$  ортодополняемо в  $M$ .

Здесь, как обычно, через  $\text{Ker } A$  обозначим ядро оператора  $A$ , а через  $A|M$  — сужение оператора  $A$  на подпространство  $M$ .

Доказательство тривиально.

**Замечание 2.11.** Если  $\lambda$  — собственное значение самосопряженного оператора  $A$ , то из ортодополняемости  $\text{Ker } (A - \lambda)$  следует, что  $\lambda$  является вещественным и полупростым собственным значением оператора  $A$ . Заметим также, что в конечномерном пространстве  $H$ , подпространство  $\text{Ker } (A - \lambda)$  ортодополняемо тогда и только тогда, когда  $\lambda$  является вещественным и полупростым собственным значением оператора  $A$ .

**Теорема 2.12.** Пусть  $A$  — самосопряженный оператор в пространстве Понтрягина  $H$ . Тогда оператор  $A$  фундаментально приводим в том и только том случае, если подпространство  $\text{Ker } (A - \lambda)$  ортодополняемо при любом собственном значении  $\lambda$ .

**Доказательство.** Достаточно доказать, что при выполнении условия теоремы, оператор  $A$  фундаментально приводим. Будем считать, что  $k = k^+(H)$  (см. стр. 51 из [8]). Заметим, что условие обеспечивает вещественность всех собственных значений. Пусть  $\lambda_1, \dots, \lambda_r$  — множество всех собственных значений оператора  $A$ , имеющих хотя бы один неотрицательный собственный вектор.

По условию мы имеем  $H = \text{Ker } (A - \lambda_1) (\dot{+}) (\text{Ker } (A - \lambda_1))^\perp$ . Обозначим подпространство  $(\text{Ker } (A - \lambda_1))^\perp$  через  $M_1$ . Подпространство  $M_1$  инвариантно относительно оператора  $A - \lambda_2$ . По лемме 2.10

$$M_1 = \text{Ker } ((A - \lambda_2)|M_1) (\dot{+}) (\text{Ker } (A - \lambda_2)|M_1)^\perp \cap M_1.$$

Обозначим подпространство  $(\text{Ker } (A - \lambda_2)|M_1)^\perp$  через  $M_2$ . Тогда точно таким же образом проверяем, что  $M_2$  инвариантно относительно оператора  $A - \lambda_3$ , и

$$M_2 = \text{Ker } ((A - \lambda_3)|M_2) (\dot{+}) (\text{Ker } (A - \lambda_3)|M_2)^\perp \cap M_2.$$

Далее, повторяя этот процесс мы получим после конечного числа шагов разложение пространства  $H$  на некоторые инвариантные подпространства оператора  $A$ :

$$H = \text{Ker } (A - \lambda_1) (\dot{+}) \text{Ker } ((A - \lambda_2)|M_1) (\dot{+}) \text{Ker } ((A - \lambda_3)|M_2) (\dot{+}) \dots \\ (\dot{+}) \text{Ker } ((A - \lambda_r)|M_{r-1}) (\dot{+}) (\text{Ker } (A - \lambda_r)|M_{r-1})^\perp \cap M_{r-1}.$$



Если бы последняя компонента этого разложения была бы индефинитной, то по теореме Понтрягина об инвариантных максимальных семидефинитных подпространствах существовал бы такой вектор  $x$  этой компоненты, что  $(x, x) \cong 0$ ,  $x \neq 0$ , и  $Ax = \lambda x$ . Отсюда  $\lambda = \lambda_i$  для некоторого  $i$ .

Но, с другой стороны, из построения разложения ясно, что  $x \in M_j$  при всех  $j$ , т.е.  $\lambda \neq \lambda_j$  при всех  $j$ . Тем самым мы пришли к противоречию.

Таким образом, последняя компонента разложения должна быть отрицательно определенной. Остальные компоненты мы разложим на пару ортогональных определенных подпространств, каждое из которых, очевидно, инвариантно относительно оператора  $A$ . Теорема доказана.

Из этой теоремы вытекает, что в конечномерном пространстве самосопряженный оператор фундаментально приводим тогда и только тогда, когда все его собственные значения вещественны и полупросты.

*Замечание 2.13. Последняя теорема верна для нормального и изометрического оператора.*

*Замечание 2.14. Компактный самосопряженный оператор в пространстве Понтрягина фундаментально приводим в том и только том случае, если его ядро ортогополняемо и все собственные значения вещественны и полупросты.*

Теперь мы приходим к вопросу, который связан с фундаментальной приводимостью оператора.

Наследуется ли фундаментальная приводимость? Точнее, пусть оператор  $T$  фундаментально приводим, и  $M$ —ортогополняемое подпространство, инвариантное относительно оператора  $T$ . Следует ли отсюда, что  $T|M$  фундаментально приводим?

*Теорема 2.15. Для самосопряженного оператора  $A$  фундаментальная приводимость наследуется.*

*Доказательство.* Пусть  $M$ —ортогополняемое инвариантное подпространство оператора  $A$ . Поскольку  $A$  фундаментально приводим, то спектр оператора  $A|M$  вещественен. По теореме 2.8, в некоторой  $J$ -норме имеет место неравенство

$$\|(A - i\gamma)^{-k}\|_J \leq \frac{C}{|\gamma|^k} \quad (k = 1, 2, \dots, \gamma \neq 0, \gamma \in \mathbb{R}).$$

Далее, легко видно, что  $(A - z)^{-1}M \subset M$  при  $\text{Im } z \neq 0$ , и вышеприведенное неравенство имеет место в  $M$  для некоторой фундаментальной симметрии пространства  $M$ . По той же теореме 2.8 оператор  $A|M$  фундаментально приводим.

*Теорема 2.16. Пусть  $U$ —фундаментально приводимый унитарный оператор в  $H$ , спектр которого является собственным подмножеством единичной окружности. Тогда для  $U$  фундаментальная приводимость наследуется.*

*Доказательство.* Поскольку оператор  $U$  фундаментально приводим, то  $UJ = JU$ , где  $J$ —некоторая фундаментальная симметрия. Отсюда, с одной стороны, мы имеем  $U^{-1} = JU^{-1}J$ , а с другой стороны, в силу унитарности,  $U^+ =$



$=JU^{-1}J$ . Итак,  $U$  является унитарным оператором в  $J$ -скалярном произведении. В частности, спектр оператора  $U$  лежит на единичной окружности.

Пусть  $M$ —ортодополняемое подпространство, инвариантное относительно оператора  $U$ . Тогда  $UJM = JUM \subset JM$ .

По теореме Уэрмера (см. [11], стр. 86), подпространство  $JM$  также инвариантно относительно оператора  $U^+$ . Отсюда, легко видно, что  $U^{-1}M \subset M$ , т.е.  $U$  является унитарным оператором в  $M$ .

По теореме 1.5 оператор  $U|_M$  фундаментально приводим.

Рассмотрим теперь следующий вопрос: когда для семейства фундаментально приводимых операторов существует фундаментальное разложение, приводящее каждый оператор из этого семейства?

**Теорема 2.17.** Пусть  $\{U_i\}_{i=1}^n$ —конечное коммутативное семейство фундаментально приводимых унитарных операторов. Тогда существует фундаментальное разложение, приводящее их одновременно.

**Доказательство.** Образует группу с образующими

$$U_1, \dots, U_n; \quad U_1^{-1}, \dots, U_n^{-1}, I.$$

Эта группа коммутативна. Действительно, оператор  $U_j$  является унитарным оператором в некотором  $J$ -скалярном произведении; но тогда по теореме Фуглида (см. [12], стр. 327) из  $U_i U_j = U_j U_i$  следует  $U_i U_j^{-1} = U_j^{-1} U_i$ .

Построенная группа ограничена по норме по теореме 1.5. По теореме 1.6 доказательство завершается.

С помощью преобразования Кэли можно доказать аналогичную теорему для самосопряженного оператора.

**Пример 2.18.** Если  $J_1$  и  $J_2$ —две различные фундаментальные симметрии, тогда  $J_1 J_2 \neq J_2 J_1$  и для  $J_1, J_2$  нет фундаментального разложения, приводящее их одновременно.

**Пример 2.19.** Этот пример показывает, что для бесконечного семейства операторов аналог последней теоремы, вообще говоря, не верен.

Рассмотрим положительный оператор, построенный в примере 4 из [1]. Если  $E(t)$  его спектральная функция, то значения этой функции образуют коммутативное семейство фундаментально приводимых положительных операторов. Если бы существовало фундаментальное разложение, приводящее все операторы из этого семейства, то из спектрального представления этого положительного оператора следовало бы, что этот оператор фундаментально приводим. Но мы знаем, что он не является фундаментально приводимым.

## ЛИТЕРАТУРА

- [1] Баясгалан, Ц., О фундаментальной приводимости положительных операторов в пространствах с индефинитной метрикой, *Studia Sci. Math. Hungar.* 13 (1978), 143—150.
- [2] Баясгалан, Ц., *Кандидатская диссертация*, Будапешт, 1980.
- [3] PHILLIPS, R. S., The extension of dual subspaces invariant under an algebra, *Proc. Internat. Sympos. Linear Spaces*, Jerusalem and Oxford, Jerusalem Academic Press and Pergamon, 1961, 366—398. MR 24#A3512.

- [4] PESONEN, E., Über die Spektraldarstellung quadratischer Formen in linearen Räumen mit indefiniter Metrik, *Ann. Acad. Sci. Fenn. Ser. A. I.* no. 227 (1956), 1—31. *MR* 20 # 1202.
- [5] KÜHNE, R., Über eine Klasse  $J$ -selbstadjungierter Operatoren, *Math. Ann.* 154 (1964), 56—69. *MR* 30 # 449.
- [6] HESS, P., Zur Theorie der linearen Operatoren eines  $J$ -Raumes. Operatoren die von kanonischen Zerlegungen reduziert werden, *Math. Z.* 106 (1968), 88—96. *MR* 38 # 567.
- [7] АКОПЯН, Р. В., К теории спектральной функции  $J$ -неотрицательного оператора, *Izv. Akad. Nauk Armjan. SSR. Ser. Mat.* 13 (1978), 114—121. *MR* 80a: 47055.
- [8] BOGNÁR, J., *Indefinite inner product spaces*, Springer-Verlag, Berlin—Heidelberg—New York, 1974. *MR* 57 # 7125.
- [9] YOSIDA, K., *Functional analysis*, Springer-Verlag, New York, 1968. *MR* 39 # 741.
- [10] SZ.-NAGY, B., On uniformly bounded linear transformations in Hilbert space, *Acta Univ. Szeged. Sect. Sci. Math.* 11 (1947), 152—157. *MR* 9—191.
- [11] DUNFORD, N. and SCHWARTZ, J. T., *Linear operators, Part II*, J. Wiley, New York, 1963. *MR* 32 # 6181.
- [12] RUDIN, W., *Functional analysis*, McGraw-Hill, New York, 1973. *MR* 51 # 1315.
- [13] MCENNIS, B. W., Fundamental reducibility of selfadjoint operators on Kreĭn space, *J. Operator Theory* 8 (1982), 219—225.

(Поступила 20-ого сентября 1983 г.)

КАФЕДРА МАТЕМАТИЧЕСКОГО АНАЛИЗА  
МОНГОЛЬСКОГО УНИВЕРСИТЕТА  
УЛАН БАТОР  
MONGOLIA



# ARCHIMEDEAN DECOMPOSITIONS OF LEFT S-SEMIMODULES AND SEMIRINGS

FRANCISCO POYATOS

A *semiring* is defined to be an algebra  $S=(S, +, \cdot)$  such that  $(S, +)$  and  $(S, \cdot)$  are semigroups and both distributive laws,  $r(s+t)=rs+rt$  and  $(s+t)r= sr+tr$  are satisfied (cf. [1]). In this note, all semirings are assumed to have commutative addition.

An additively written commutative semigroup  $M=(M, +)$  is called a *semimodule*. If, for each element  $s$  of a semiring  $S$ , there is an unary operation  $f_s: M \rightarrow M$  denoted by  $f_s(a)=sa$  such that

$$s(a+b) = sa+sb, \quad (s+t)a = sa+ta, \quad (st)a = s(ta)$$

hold for all  $s, t \in S; a, b \in M$ , then  ${}_sM=(M, +, \{f_s\}_{s \in S})$  is called a *left semimodule over  $S$* , or a *left  $S$ -semimodule* (cf. [3], [4]). With respect to these algebras we speak about  $S$ -subsemimodules,  $S$ -homomorphisms, etc. In particular, an  $S$ -congruence  $\kappa$  on a left  $S$ -semimodule  ${}_sM$  is a congruence on  $M=(M, +)$  which also satisfies

$$(1) \quad axb \Rightarrow saxsb \quad \text{for all } a, b \in M; s \in S.$$

Moreover, a left  $S$ -semimodule  ${}_sM$  is called a *left  $S$ -semilattice*, if  $(M, +)$  is idempotent. Right  $S$ -semimodules  $M_S$  are defined dually.

The purpose of this note is to generalize a well-known structure theorem on commutative semigroups, due to T. Tamura, N. Kimura and G. Thierrin (cf. [2], Theorem 4.13), to  $S$ -semimodules and semirings.

**THEOREM 1.** *Let  ${}_sM=(M, +, \{f_s\}_{s \in S})$  be a left semimodule over a semiring  $S$  and*

$$(2) \quad M = \bigcup_{\alpha \in Y} M_\alpha, \quad M_\alpha + M_\beta \subseteq M_{\alpha+\beta}$$

*the unique decomposition of the semimodule  $(M, +)$  as a semilattice  $Y=(Y, +)$  of archimedean subsemimodules  $M_\alpha$  ( $\alpha \in Y$ ). Then  $Y$  is a left  $S$ -semilattice and the maximal idempotent  $S$ -homomorphic image of  ${}_sM$ , and the equivalence classes  $M_\alpha$  of the corresponding  $S$ -congruence  $\eta$  on  ${}_sM$  satisfy*

$$(3) \quad sM_\alpha \subseteq M_{sa} \quad \text{for all } s \in S.$$

1980 *Mathematics Subject Classification*. Primary 16A78; Secondary 20M14.

*Key words and phrases*. Archimedean decomposition, Tamura—Kimura—Thierrin, maximal homomorphic image.

PROOF. According to the theorem cited above, the unique decomposition (2) of  $(M, +)$  corresponds to the congruence  $\eta$  on  $(M, +)$  defined by

$$(4) \quad a\eta b \Leftrightarrow a+x=mb \quad \text{and} \quad b+y=na \quad \text{for some } x, y \in M$$

and some positive integers  $n, m$ .

Then  $Y \simeq M/\eta$  is the maximal idempotent homomorphic image of  $(M, +)$  and the equivalence classes  $M_\alpha$  of  $M$  modulo  $\eta$  are the archimedean components of  $(M, +)$ . For a left  $S$ -semimodule  ${}_S M$ , (4) is obviously also an  $S$ -congruence on  ${}_S M$  according to (1), which implies our theorem. As a consequence of (3) we note:

COROLLARY 1. *An archimedean component  $M_\alpha$  of a left  $S$ -semimodule  ${}_S M$  is an  $S$ -subsemimodule of  ${}_S M$  iff  $M_\alpha$  contains an element  $a$  such that  $sa\eta a$  holds for all  $s \in S$ , which is then true for all  $a \in M_\alpha$ .*

THEOREM 2. *Let  $S=(S, +, \cdot)$  be a semiring and*

$$(2') \quad S = \bigcup_{\alpha \in Y} S_\alpha, \quad S_\alpha + S_\beta \subseteq S_{\alpha+\beta}$$

*the unique decomposition of the semimodule  $(S, +)$  as a semilattice  $Y=(Y, +)$  of archimedean subsemimodules  $S_\alpha$  ( $\alpha \in Y$ ). Then  $Y$  is a semiring and the maximal additively idempotent semiring-homomorphic image of  $S$ , and the equivalence classes  $S_\alpha$  of the corresponding semiring-congruence  $\eta$  of  $S$  satisfy*

$$(3') \quad sS_\alpha \subseteq S_{s\alpha} \quad \text{and} \quad S_\alpha s \subseteq S_{\alpha s} \quad \text{for all } s \in S$$

*and hence  $S_\alpha \cdot S_\beta \subseteq S_{\alpha \cdot \beta}$ .*

PROOF. Each semiring  $S$  may be considered as a left  $S$ -semimodule  ${}_S S = (S, +, \{f_s\}_{s \in S})$  where  $f_s(a) = sa$  is defined by the semiring multiplication, and dually as a right  $S$ -semimodule  $S_S$ . Thus Theorem 2 follows from Theorem 1 and its dual. From the last statement we obtain:

COROLLARY 2. *An archimedean component  $S_\alpha$  of a semiring  $S$  is a subsemiring of  $S$  iff  $S_\alpha$  contains an element  $s$  such that  $s^2\eta s$  holds, which is then true for all  $s \in S_\alpha$ .*

#### REFERENCES

- [1] BOURNE, S., The Jacobson radical of a semiring, *Proc. Nat. Acad. Sci. USA* 37 (1951), 163—170. *MR* 13—7.
- [2] CLIFFORD, A. H. and PRESTON, G. B., *The algebraic theory of semigroups*, Vol. I, Amer. Math. Soc., 1961. *MR* 24#A2627.
- [3] POYATOS, F., Modulos de completion de A-semimodulos, *Rev. Mat. Hisp.-Amer.* (4) 31 (1971), 123—156. *MR* 46#7326.
- [4] STEINFELD, O., Über die Struktursätze der Semiringe, *Acta Math. Acad. Sci. Hungar.* 10 (1959), 149—155. *MR* 21#7239.

(Received October 19, 1983)

INSTITUTO "JORGE JUAN" DE MATEMATICAS  
CONSEJO SUPERIOR DE INVESTIGACIONES CIENTIFICAS  
SERRANO, 123  
E-28006 MADRID  
SPAIN

# THE LARGEST COMPONENTS IN A RANDOM LATTICE

G. R. GRIMMETT

## Abstract

Let  $p$  satisfy  $0 < p < 1$ , and let  $H$  be the random graph obtained by deleting each edge of the square lattice  $\mathbb{Z}^2$  with probability  $1-p$ . We study the sizes of the largest components of the finite subgraph  $H(n)$  of  $H$  induced by the square subsection  $B(n) = \{(i, j): 1 \leq i, j \leq n\}$  of  $\mathbb{Z}^2$ . We show that there exist positive quantities  $\alpha(p)$ ,  $\beta(p)$ ,  $\gamma(p)$ ,  $\delta(p)$ , depending on  $p$  alone, such that the sizes  $K_1(n)$  and  $K_2(n)$  of the largest and second-largest components of  $H(n)$ , respectively, satisfy (as  $n \rightarrow \infty$ )

if  $p < \frac{1}{2}$  then  $(\log n)^{-1} K_1(n) \rightarrow \alpha(p)$  in probability,

if  $p > \frac{1}{2}$  then  $n^{-2} K_1(n) \rightarrow \beta(p)$  in probability

and  $P(\gamma(p) \leq (\log n)^{-2} K_2(n) \leq \delta(p)) \rightarrow 1$ .

## 1. Introduction

Given an initial graph  $G=(V, E)$ , we may obtain a random subgraph  $H$  of  $G$  as follows. Fix  $p$  such that  $0 < p < 1$ , and delete each edge in  $E$  with probability  $1-p$ , independently of all other edges; we denote by  $H$  the resulting (random) subgraph of  $G$ , comprising the vertex set  $V$  together with all remaining edges. Of especial interest are the two cases when  $G$  is the complete graph on  $n$  vertices, or  $G$  is either the whole or part of a crystalline lattice (such as the square lattice  $\mathbb{Z}^2$ ). The former case of random subgraphs of  $K_n$  is well-studied and well understood (see Erdős and Rényi (1960) for the pioneering work, and Grimmett (1983) or Karoński (1982) for recent reviews); for this case, several celebrated results deal with the sizes of the largest and second largest components (see Erdős and Rényi (1960) and Komlós, Sulyok and Szemerédi (1980)). In this paper we consider the second case, and we assume henceforth that  $G=\mathbb{Z}^2$ , the square lattice with vertex set  $\{(i, j): i, j = 0, \pm 1, \pm 2, \dots\}$  and edges joining pairs  $(i, j)$ ,  $(k, l)$  of vertices whenever  $|i-k| + |j-l| = 1$ . The random subgraph  $H$  of  $\mathbb{Z}^2$  is usually called the *bond percolation model*, having been proposed by Broadbent and Hammersley (1957) as a model for the flow of liquid through a porous medium. Let  $B(n) = \{(i, j): 1 \leq i, j \leq n\}$  be the square subsection of  $\mathbb{Z}^2$  with side length  $n-1$  and bottom left-hand corner at  $(1, 1)$ , and let  $H(n)$  be the subgraph of  $H$  induced by  $B(n)$ . We shall study the asymptotic properties, as  $n \rightarrow \infty$ , of the sizes  $K_1(n)$  and  $K_2(n)$  of the largest and second-largest components of  $H(n)$ , and shall indicate how these properties depend

on the numerical value of the edge-probability  $p$ . It turns out that if  $p < 1/2$  then  $K_1(n)$  and  $K_2(n)$  are about  $\alpha(p) \log n$ , whilst if  $p > 1/2$  then  $K_1(n)$  is about  $\beta(p)n^2$  and  $K_2(n)$  has order  $(\log n)^2$ ; here  $\alpha(p)$  and  $\beta(p)$  are positive quantities which depend on  $p$  only. These results follow fairly simply from currently-known facts about the percolation model, and they extend results of Füredi (1979). We have no interesting results for the case  $p = 1/2$ .

We conclude the introduction by recalling some well-known properties of the percolation model on  $\mathbb{Z}^2$  which explain the dependence of  $K_1(n)$  and  $K_2(n)$  upon whether  $p < 1/2$  or  $p > 1/2$ . For any given  $p$ , we denote by  $P_p$  the corresponding probability measure. For fixed  $p$ , let  $H$  be the random subgraph of  $\mathbb{Z}^2$  as described above, and let  $W$  be the size of the component of  $H$  which contains the origin  $(0, 0)$ . We define

$$(1.1) \quad \beta(p) = P_p(W = \infty)$$

to be the probability that the origin is in an infinite component. It is the case that

$$\beta(p) \begin{cases} = 0 & \text{if } p \leq \frac{1}{2} \\ > 0 & \text{if } p > \frac{1}{2} \end{cases}$$

and  $\beta: [0, 1] \rightarrow [0, 1]$  is a continuous function (see Kesten (1980) and Russo (1978)). Furthermore, if  $p \leq 1/2$  then all components of  $H$  are almost surely (a.s.) finite, whilst if  $p > 1/2$  then  $H$  contains a.s. a unique infinite component. For recent results on the percolation model, we refer the reader to the book by Kesten (1982). It is upon the square lattice that most interest has been concentrated, and more is known about this lattice than any other. Similar results hold for certain other two-dimensional lattices, and partial results are known for higher-dimensional lattices  $\mathbb{Z}^d$ . We consider the case of  $\mathbb{Z}^2$  only here, noting that corresponding results hold similarly in the other cases whenever the necessary facts about the percolation model are known. We explore convergence in probability only.

This paper deals with the size of the largest components of  $H(n)$ . This may be contrasted with the results of Grimmett (1981b) and Kesten (1981), who studied the number of components of  $H(n)$ ; further results on this quantity have been obtained by Cox and Grimmett (1981, 1983) in a more general setting in which square boxes are replaced by the interiors of arbitrary circuits.

Related results appear in Révész (1980) and Nemetz and Kusolitsch (1982) who considered the largest red square and the largest red rectangle, respectively, in  $B(n)$  when each vertex of  $\mathbb{Z}^2$  is coloured red or blue with probability  $p$  or  $1-p$ .

## 2. The results

We prove the following two theorems. All limits are taken as  $n \rightarrow \infty$ .

**THEOREM 1.** *If  $p < 1/2$  then there exists a positive quantity  $\alpha(p)$ , depending on  $p$  only, such that*

$$(2.1) \quad \frac{1}{\log n} K_1(n) \rightarrow \alpha(p) \quad \text{in probability } (P_p).$$



THEOREM 2. If  $p > 1/2$  then

$$(2.2) \quad \frac{1}{n^2} K_1(n) \rightarrow \beta(p) \quad \text{in probability } (P_p)$$

where  $\beta$  is given by (1.1), and there exist positive quantities  $\gamma(p)$ ,  $\delta(p)$ , depending on  $p$  only, such that

$$(2.3) \quad P_p \left( \gamma(p) \leq \frac{1}{(\log n)^2} K_2(n) \leq \delta(p) \right) \rightarrow 1.$$

The following facts are straightforward consequences of the proofs of these theorems.

(a) Suppose  $p < 1/2$ . For any  $\varepsilon > 0$ , the number of components of  $H(n)$  with sizes in  $[(1-\varepsilon)\alpha(p)\log n, (1+\varepsilon)\alpha(p)\log n]$  tends to infinity as  $n \rightarrow \infty$ , with probability  $1 - o(1)$ .

(b) Suppose  $p > 1/2$ . If  $\gamma(p)$  and  $\delta(p)$  are chosen correctly in (2.3) then the number of components of  $H(n)$  with sizes satisfying the inequalities of (2.3) tends to infinity as  $n \rightarrow \infty$ , with probability  $1 - o(1)$ .

We have presented Theorems 1 and 2 in their simplest forms. Some minor improvements may be made to them by estimating the rates of convergence in (2.1) and (2.2). These theorems extend results of Füredi (1979), who stated the following facts. There exist positive quantities  $\mu(p)$ ,  $\sigma(p)$  such that

$$(a) \quad \text{if } p < \frac{1}{3} \quad \text{then } P_p(K_1(n) < \mu(p)\log n) \rightarrow 1,$$

$$(b) \quad \text{if } p > 2/3 \quad \text{then } H(n) \text{ contains a giant component and}$$

$$P_p(K_2(n) < \sigma(p)(\log n)^2) \rightarrow 1.$$

### 3. The proofs

Let

$$(3.1) \quad \pi_p(n) = P_p(|W| = n), \quad \Pi_p(n) = P_p(n \leq |W| < \infty).$$

Of course,  $\Pi_p(n) = P_p(|W| \geq n)$  if  $p \leq 1/2$ . We shall use the following results from percolation theory. Kunz and Souillard (1978, p. 91) proved that there exists  $A(p)$  satisfying  $0 < A(p) \leq \infty$  such that

$$(3.2) \quad -\frac{1}{n} \log \pi_p(n) \rightarrow A(p).$$

As a consequence of Theorem 1 of Kesten (1981), we have that

$$(3.3) \quad \text{if } p < \frac{1}{2} \quad \text{then } 0 < A(p) < \infty.$$

If  $p > 1/2$  then  $A(p) = \infty$ ; in this case, the asymptotic behaviour of the sequences  $\{\pi_p(n)\}$  and  $\{\Pi_p(n)\}$  are such that there exist positive quantities  $\Gamma(p)$ ,  $\Delta(p)$ , depend-

ing on  $p$  alone, such that, for all large  $n$ ,

$$(3.4) \quad \Delta(p) \equiv -\frac{1}{\sqrt[n]{n}} \log \Pi_p(n) \equiv -\frac{1}{\sqrt[n]{n}} \log \pi_p(n) \equiv \Gamma(p);$$

see Aizenman, Delyon and Souillard (1980) and Kesten (1982, pp. 98, 99) for these inequalities. Less is known in the case when  $p=1/2$ , and for that reason we do not consider that case here; see Kesten (1982, p. 227) for a brief discussion.

In the proofs which follow, we shall often use non-integer-valued quantities in places where integers are required; it will be clear that this makes no essential difference to the argument.

**PROOF OF THEOREM 1.** Assume that  $p < 1/2$ . Define

$$(3.5) \quad \alpha(p) = 2A(p)^{-1},$$

and suppose  $\varepsilon$  satisfies  $0 < \varepsilon < 1/2$ . Then, for all large  $n$ ,

$$\begin{aligned} P_p(K_1(n) > (1+2\varepsilon)\alpha(p)\log n) &\leq n^2 \Pi_p((1+2\varepsilon)\alpha(p)\log n) \leq \\ &\leq Cn^{2-(1-\varepsilon)(1+2\varepsilon)A(p)\alpha(p)} = \\ &= Cn^{-2\varepsilon(1-2\varepsilon)} \rightarrow \\ &\rightarrow 0, \end{aligned}$$

by (3.2) and (3.5), where  $C=C(p)$  is a constant. To see that  $K_1(n)$  cannot be too small, we partition  $B(n)$  into subsquares with side lengths  $2\alpha(p)\log n$ ; there are

$$\left( \frac{n}{2\alpha(p)\log n} \right)^2 = N$$

of these subsquares, which we denote by  $B_1, B_2, \dots, B_N$ . For each  $i$ , let  $x_i$  be a vertex which is as close as possible to the centre of  $B_i$ , and let  $I_i$  be the indicator function of the event that  $x_i$  is in a component of  $H(n)$  with size at least  $(1-\varepsilon)\alpha(p)\log n$ . For all large  $n$ , the random variable  $I_i$  depends on the presence or absence of the edges in  $B_i$  only, and thus  $I_1, I_2, \dots, I_N$  are independent. It follows that the sum

$$S_N = I_1 + I_2 + \dots + I_N$$

is binomially distributed with parameters  $N$  and  $\Pi_p((1-\varepsilon)\alpha(p)\log n)$ . Therefore, for all large  $n$ ,

$$\begin{aligned} (3.6) \quad P_p(K_1(n) < (1-\varepsilon)\alpha(p)\log n) &\leq P(S_N = 0) = \\ &= (1 - \Pi_p((1-\varepsilon)\alpha(p)\log n))^N \leq \\ &\leq (1 - n^{-(1-\varepsilon^2)A(p)\alpha(p)})^N \leq \\ &\leq \exp \left( - \left( \frac{n^{\varepsilon^2}}{2\alpha(p)\log n} \right)^2 \right) \rightarrow \\ &\rightarrow 0 \end{aligned}$$

by (3.2) and (3.5). The proof is complete.

PROOF OF THEOREM 2. Assume that  $p > 1/2$ , and write  $B(m, n)$  for the rectangle  $\{(i, j): 1 \leq i \leq m, 1 \leq j \leq n\}$  and  $H(m, n)$  for the subgraph of  $H$  induced by  $B(m, n)$ . We write  $\{B(m, n) \text{ horizontally}\}$  (respectively  $\{B(m, n) \text{ vertically}\}$ ) for the event that  $B(m, n)$  is traversed horizontally (respectively vertically) by a path of  $H(m, n)$  joining two vertices on the two corresponding opposite sides of  $B(m, n)$ . From Theorem 2 of Grimmett (1981a), there exists  $v(p) > 0$  such that

$$(3.7) \quad P_p(B(v(p) \log n, n) \text{ vertically}) = P_p(B(n, v(p) \log n) \text{ horizontally}) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

Henceforth we suppose that  $v(p)$  is chosen such that (3.7) holds. Let  $C(n)$  be the square

$$C(n) = \{(i, j): v(p) \log n < i, j < n - v(p) \log n\}$$

and let  $A(n) = B(n) \setminus C(n)$  denote the annulus of  $B(n)$  which surrounds  $C(n)$ .  $A(n)$  may be written as the union  $A(n) = R_1 \cup R_2 \cup R_3 \cup R_4$  of four overlapping rectangles, each with dimensions  $n$  by  $v(p) \log n$ . Let  $\mathcal{A} = \{A_1, A_2, \dots\}$  be the collection of all circuits of  $H$  which lie in  $A(n)$  and have the property that they contain  $C(n)$  in their interiors. By (3.7) and the FKG inequality (see Kesten (1982, p. 72))

$$(3.8) \quad P_p(\mathcal{A} \text{ is empty}) \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

If  $\mathcal{A}$  is non-empty, let  $K$  be the set of vertices of  $B(n)$  which either lie outside all the circuits in  $\mathcal{A}$  but are joined to some  $A_i$  by a path in  $H(n)$ , or which lie inside some  $A_i$  and are in the infinite component of  $H$ . By the planarity of  $\mathbb{Z}^2$ ,  $K$  is the vertex set of a component of  $H(n)$ . Next we estimate the size of  $K$ . For each  $(i, j) \in \mathbb{Z}^2$  we define  $I_{ij}$  to be the indicator function of the event that  $(i, j)$  is in the infinite component of  $H$ . The collection  $\{I_{ij}: (i, j) \in \mathbb{Z}^2\}$  is a stationary family of random variables, and it follows by a suitable ergodic theorem (see, for example, Dunford (1951)) that the number

$$K(n) = \sum_{(i, j) \in C(n)} I_{ij}$$

of vertices in  $C(n)$  which belong to the infinite component of  $H$  satisfies

$$(3.9) \quad \frac{1}{|C(n)|} K(n) \rightarrow E_p(I_{00}) = \beta(p) \quad \text{in probability,}$$

where

$$\frac{1}{n^2} |C(n)| \rightarrow 1.$$

From (3.8), (3.9), and the observations in between, it follows that, for all  $\varepsilon > 0$ ,

$$P_p(K_1(n) > (1 - \varepsilon)\beta(p)n^2) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

The corresponding upper bound is obvious. By (3.8) and (3.9), any component with size at least  $(1 + \varepsilon)\beta(p)n^2$  in  $H(n)$ , where  $\varepsilon > 0$ , contains (with probability  $1 - o(1)$ ) a vertex in  $C(n)$  which is in a finite component of  $H$ . The probability that such a

vertex exists is no bigger than

$$n^2 \Pi_p((1+\varepsilon)\beta(p)n^2) \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

by (3.4).

The proof of (2.3) is similar to the proof of (2.1), using (3.4) in place of (3.2). The first inequality of (2.3) is valid so long as  $\gamma(p)$  satisfies

$$\gamma(p) < 4\Gamma(p)^{-2},$$

and the proof of this differs from that of (3.6) only in that we divide  $B(n)$  into subsquares with side-length  $(\log n)^3$ .

To see the other part, let  $\Delta B(n)$  be the boundary of  $B(n)$ , defined by  $\Delta B(n) = \{x \in B(n) : x \text{ is adjacent in } \mathbb{Z}^2 \text{ to some } y \notin B(n)\}$ .

The boundary  $\Delta C(n)$  of  $C(n)$  is defined similarly. If  $K_2(n) > \delta(p)(\log n)^2$  and  $\mathcal{A}$  is non-empty then  $H(n)$  has a component  $K'$  (other than  $K$ ) of size at least  $\delta(p)(\log n)^2$ , and either  $K'$  intersects  $\Delta B(n)$  or  $K'$  does not intersect  $\Delta B(n)$ . If  $K'$  does not intersect  $\Delta B(n)$  then some vertex of  $B(n)$  is in a finite component of  $H$  with size exceeding  $\delta(p)(\log n)^2$ , and the probability of this is no larger than

$$n^2 \Pi_p(\delta(p)(\log n)^2) \leq n^{2-\Delta(p)\sqrt{\delta(p)}} \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

if

$$(3.10) \quad \delta(p) > 4\Delta(p)^{-2},$$

by (3.4). On the other hand, if  $K'$  intersects  $\Delta B(n)$  then  $K'$  is contained strictly within  $A(n)$  (since  $\mathcal{A}$  is assumed to be non-empty); we use the dual-lattice technique (see Kesten 1982, p. 37) to see that this is unlikely. There exists an absolute constant  $\varrho$  with the property that, given any component of  $H(n)$  in  $A(n)$  with size at least  $\delta(p)(\log n)^2$ , there exists a path in the dual graph of  $A(n) \cup \Delta C(n)$  which has length at least  $\varrho(\delta(p)(\log n)^2)^{1/2} = \varrho\delta(p)^{1/2} \log n$  and which crosses no edge of  $\mathbb{Z}^2$  which lies in  $H$ . But the dual of  $\mathbb{Z}^2$  is isomorphic to  $\mathbb{Z}^2$ , and so the probability of this is no larger than

$$\begin{aligned} |A(n) \cup \Delta C(n)| \Pi_{1-p}(\varrho\delta(p)^{1/2} \log n) &\leq \\ &\leq 4v(p) n \log n \exp(-(1-\varepsilon)A(1-p)\varrho\delta(p)^{1/2} \log n) \end{aligned}$$

for all large  $n$ , where  $\varepsilon > 0$ . This probability tends to 0 as  $n \rightarrow \infty$  if

$$\delta(p) > ((1-\varepsilon)\varrho A(1-p))^{-2},$$

giving from (3.10) that (2.3) holds whenever

$$\delta(p) > \max \{4\Delta(p)^{-2}, 2(\varrho A(1-p))^{-2}\}.$$

## REFERENCES

- [1] AIZENMAN, M., DELYON, F. and SOUILLARD, B., Lower bounds on the cluster size distribution, *J. Statist. Phys.* **23** (1980), 267–280. *MR* **82b**: 82048.
- [2] BROADBENT, S. R. and HAMMERSLEY, J. M., Percolation processes. I. Crystals and mazes, *Proc. Cambridge Philos. Soc.* **53** (1957), 629–641. *MR* **19**: 989.
- [3] COX, J. T. and GRIMMETT, G. R., Central limit theorems for percolation models, *J. Statist. Phys.* **25** (1981), 237–251. *MR* **82k**: 60209.

- [4] COX, J. T. and GRIMMETT, G. R., Central limit theorems for associated random variables and the percolation model, *Ann. Probability* **12**(1984), 514—528.
- [5] DUNFORD, N., An individual ergodic theorem for noncommutative transformations, *Acta Sci. Math. Szeged* **14** (1951), 1—5. *MR* **13**—49.
- [6] ERDŐS, P. and RÉNYI, A., On the evolution of random graphs, *Magyar Tud. Akad. Mat. Kutató Int. Közl.* **5** (1960), 17—61. *MR* **23**#A2338.
- [7] FÜREDI, Z., On connectedness of a random graph with small number of edges, *Studia Sci. Math. Hungar.* **14** (1979), 419—425.
- [8] GRIMMETT, G. R., Critical sponge dimensions in percolation theory, *Adv. in Appl. Probab.* **13** (1981), 314—324. *MR* **82k**: 60210.
- [9] GRIMMETT, G. R., On the differentiability of the number of clusters per vertex in the percolation model, *J. London Math. Soc.* **23** (1981), 372—384. *MR* **82g**: 60135.
- [10] GRIMMETT, G. R., Random graphs, in *Selected topics in graph theory II*, ed. L. Beineke and R. Wilson, Academic Press, London, 1983, 201—235.
- [11] KARONSKI, M., A review of random graphs, *J. Graph Theory* **6** (1982), 349—389.
- [12] KESTEN, H., The critical probability of bond percolation on the square lattice equals  $1/2$ , *Comm. Math. Phys.* **74** (1980), 41—59. *MR* **82c**: 60179.
- [13] KESTEN, H., Analyticity properties and power law estimates of functions in percolation theory, *J. Statist. Phys.* **25** (1981), 717—756. *MR* **83b**: 82062.
- [14] KESTEN, H., *Percolation theory for mathematicians*, Birkhäuser, Boston, 1982.
- [15] KOMLÓS, J., SÜLYÖK, M. and SZEMERÉDI, E., Second largest component in a random graph, *Studia Sci. Math. Hungar.* **15** (1980), 391—395.
- [16] KUNZ, H. and SOUILLARD, B., Essential singularity in percolation problems and asymptotic behaviour of cluster size distribution, *J. Statist. Phys.* **19** (1978), 77—106. *MR* **58**#14852.
- [17] NEMETZ, T. and KUSOLITSCH, N., On the longest run of coincidences, *Z. Wahrscheinlichkeitstheorie und Verw. Gebiete* **61** (1982), 59—73.
- [18] RÉVÉSZ, P., How to characterize the asymptotic properties of a stochastic process by four classes of deterministic curves?, preprint, (1980).
- [19] RUSSO, L., A note on percolation, *Z. Wahrscheinlichkeitstheorie und Verw. Gebiete* **43** (1978), 39—48. *MR* **58**#7931.

( Received October 24, 1983 )

SCHOOL OF MATHEMATICS  
UNIVERSITY OF BRISTOL  
UNIVERSITY WALK  
BRISTOL BS8 1TW  
ENGLAND



# NORM FORM EQUATIONS WITH SEVERAL DOMINATING VARIABLES AND EXPLICIT LOWER BOUNDS FOR INHOMOGENEOUS LINEAR FORMS WITH ALGEBRAIC COEFFICIENTS II

I. GAÁL

## 1. Inhomogeneous norm form equations

In the first part of this paper (see [5]) we considered some inhomogeneous norm form equations. The purpose of the present paper is to extend the results of [5] to the case when the coefficients in the norm form satisfy certain weaker assumptions than in [5].

Let  $L, K$  be algebraic number fields with  $L \subset K$  and  $[K:L] = n \geq 3$ . Let  $\alpha_1 = 1, \alpha_2, \dots, \alpha_k$  ( $k \geq 2$ ) be elements of  $K$ , linearly independent over  $L$ , such that  $K = L(\alpha_2, \dots, \alpha_k)$ . Further let  $0 \neq \mu \in L$ . Let us consider the norm form equation

$$(1) \quad N_{K|L}(x_1 + \alpha_2 x_2 + \dots + \alpha_k x_k) = \mu$$

where the variables are<sup>1</sup>  $x_1, \dots, x_k \in \mathbb{Z}_L$ . Under the condition

$$(2) \quad [L(\alpha_i): L] = n_i \geq 3 \quad (i = 2, \dots, k), \quad \text{and} \quad n_2, \dots, n_k = n$$

Győry and Papp [16], [17] gave effective upper bounds for all solutions of equation (1). Győry and Papp [16] (see also Győry [11]) derived effective bounds for the solutions of (1) also under the weaker condition

$$(3) \quad \alpha_i \text{ is of degree } \geq 3 \text{ over } L(\alpha_1, \dots, \alpha_{i-1}), \quad i = 2, \dots, k.$$

Later Kotov [19] gave another proof for a slightly weaker version of this effective result. For certain  $p$ -adic generalizations see Győry [6], [9], [10], [12], [13] and Kotov [20]. Under the condition  $\alpha_1, \dots, \alpha_{k-1}$  are linearly independent over  $L$  and  $\alpha_k$  is of degree  $\geq 3$  over  $L(\alpha_2, \dots, \alpha_{k-1})$ , recently Győry [12], [13] and Kotov [21], [22], independently gave effective bounds for all solutions of equation (1), satisfying  $x_k \neq 0$ . Moreover, in [12], [13] and [21] the results mentioned have been generalized to include the  $p$ -adic case, too. The above quoted results of Győry and Papp, Győry and Kotov generalize and improve, respectively, many earlier results on Thue equations and Thue—Mahler equations. (See e.g. the theorems of Thue [35], Mahler [25], Baker [1], [2], Coates [3], Feldman [4], Sprindžuk [29], [31], Stark [34], Kotov [18], Kotov and Sprindžuk [23], Shorey, van der Poorten, Tijdeman and Schinzel [28].) We remark, that very recently Győry [14], [15] extended these results to the more general case, when the ground ring is an arbitrary integral domain, finitely generated over  $\mathbb{Z}$ .

<sup>1</sup>  $\mathbb{Z}_L$  denotes the ring of integers of an algebraic number field  $L$ .

1980 *Mathematics Subject Classification*. Primary 10B16, 10B45; Secondary 10F25 10F30.

*Key words and phrases*. Norm form equations, Diophantine inequalities, approximation to algebraic numbers, approximation by numbers from a fixed field.



As a common generalization of a theorem of Sprindžuk [32] (see also [33]) and the above mentioned theorem of Győry and Papp [17], in [5] we considered the inhomogeneous norm form equation

$$(4) \quad N_{K|L}(x_1 + \alpha_2 x_2 + \dots + \alpha_k x_k + \lambda) = \mu$$

where  $\alpha_2, \dots, \alpha_k$  satisfy (2). In equation (4) the dominating variables are  $x_1, \dots, x_k \in \mathbf{Z}_L$  and  $\lambda \in \mathbf{Z}_K$  is a non-dominating variable satisfying<sup>2</sup>  $|\lambda| < (\max_{2 \leq i \leq k} |\alpha_i|)^{1-\zeta}$  where  $0 < \zeta < 1$  is a given small positive constant. Under the above assumptions we obtained effective upper bounds for all solutions of equation (4). In the special case  $k=2$ ,  $L=\mathbf{Q}$  our theorem gave Sprindžuk's result [32] which is an inhomogeneous generalization of Baker's famous theorem [1] on Thue's equation. Further, in the case  $\lambda=0$  our theorem implied the result of Győry and Papp [17] which is a generalization of the above mentioned theorem of Baker to the case of norm form equations in several variables.

In this paper we give effective bounds for all solutions of equation (4) in the case  $L=\mathbf{Q}$ , assuming that  $\alpha_2, \dots, \alpha_k$  satisfy the weaker condition (3) instead of (2).

Let  $K$  be an algebraic number field of degree  $n$ . Denote by  $r$  and  $R_K$  the number of fundamental units and the regulator of  $K$  and let  $R_K^* = \max(R_K, e)$ . Let  $\alpha_1 = 1, \alpha_2, \dots, \alpha_k$  ( $k \geq 2$ ) be elements of  $K$  such that  $K = \mathbf{Q}(\alpha_2, \dots, \alpha_k)$  and<sup>3</sup>  $H(\alpha_i) \leq H$  ( $H \geq e$ ). Suppose that  $\alpha_i$  is of degree  $\geq 3$  over  $\mathbf{Q}(\alpha_1, \dots, \alpha_{i-1})$  for  $i=2, \dots, k$ . Let  $m$  be a non-zero rational number. Let us consider the equation

$$(5) \quad N_{K|\mathbf{Q}}(x_1 + \alpha_2 x_2 + \dots + \alpha_k x_k + \lambda) = m$$

where the dominating variables are  $x_1, \dots, x_k \in \mathbf{Z}$ , and  $\lambda \in \mathbf{Z}_K$  is a non-dominating variable satisfying  $|\lambda| < C(\max_{1 \leq i \leq k} |x_i|)^{1-\zeta}$  ( $C \geq e$ ,  $0 < \zeta < 1$  are given constants).

Under the above assumptions our main result is the following:

**THEOREM 1.** *If  $x_1, \dots, x_k \in \mathbf{Z}$ ,  $\lambda \in \mathbf{Z}_K$  are solutions of (5) and  $|\lambda| < C(\max_{1 \leq i \leq k} |x_i|)^{1-\zeta}$  then*

$$(6) \quad \max_{1 \leq i \leq k} |x_i| < \exp \left[ \frac{9^{k-2}}{\zeta^2} c_1 \log(H^{2n^2} |m| C) \right]$$

with

$$c_1 = (25(r+3)n)^{18(r+3)} (R_K^* \log R_K^*)^2.$$

In the special case  $\lambda=0$  our theorem includes e.g. the case  $L=\mathbf{Q}$  of a result of Győry and Papp [16] (see also Theorem 3.3 in Győry [11]) and Theorem 2 of Kotov [19].

The following corollary is an immediate consequence of our Theorem 1.

**COROLLARY 1.1.** *Suppose that in equation (5)  $[\mathbf{Q}(\alpha_i):\mathbf{Q}] = n_i \geq 3$  for  $i=2, \dots, k$  and  $n_2 \dots n_k = n$ . If  $x_1, \dots, x_k \in \mathbf{Z}$ ,  $\lambda \in \mathbf{Z}_K$  are solutions of (5) and  $|\lambda| < C(\max_{1 \leq i \leq k} |x_i|)^{1-\zeta}$  then  $x_1 \dots x_k$  satisfies (6).*

<sup>2</sup> For an algebraic number  $\alpha$ ,  $|\alpha|$  denotes the size of  $\alpha$ , that is the maximum absolute value of its conjugates.

<sup>3</sup>  $H(\alpha)$  denotes the height of an algebraic number  $\alpha$ , that is the maximum absolute value of the coefficients of the minimal defining polynomial of  $\alpha$  over  $\mathbf{Z}$ .

The statement of this corollary essentially coincides with the case  $L=\mathbf{Q}$  of Theorem 1 of [5]. In the special case  $k=2$  Corollary 1.1 gives Theorem 1 of Sprindžuk [32]. Further, in the special case  $\lambda=0$  our above corollary implies the case  $L=\mathbf{Q}$  of Theorem 1 of Györy and Papp [17].

## 2. Inhomogeneous linear forms with algebraic coefficients

As in the first part of our paper (cf. [5]), Theorem 1 makes possible to derive effective lower bounds for certain inhomogeneous linear forms with algebraic coefficients.

Let  $\alpha_1=1, \alpha_2, \dots, \alpha_k$  ( $k \geq 2$ ) be algebraic numbers, linearly independent over  $\mathbf{Q}$ . Let  $K=\mathbf{Q}(\alpha_2, \dots, \alpha_k)$  and  $[K:\mathbf{Q}]=n$ . By a generalization of a well-known theorem of Liouville [24], there exists an effectively computable positive constant  $c>0$ , depending on  $\alpha_2, \dots, \alpha_k$  such that

$$|x_1 + x_2\alpha_2 + \dots + x_k\alpha_k| > cX^{-(n-\sigma)/\sigma}, \quad X = \max_{1 \leq i \leq k} |x_i|$$

for any  $(x_1, \dots, x_k) \in \mathbf{Z}^k \setminus \{O\}$ , where  $\sigma=1$  or  $2$  according as  $K$  is real or not. Effective bounds obtained for the solutions of norm form equations make possible to give effective improvements in the exponent of the above Liouville inequality. Theorems of this type were obtained e.g. by Baker [1], Feldman [4], Sprindžuk [30], Kotov and Sprindžuk [23] in the case  $k=2$ , and, more generally, by Györy and Papp [17], Györy [8] and Kotov [21] in the case  $k \geq 2$ . In [8], [21] the authors established their estimates for linear forms with coefficients satisfying (3) and extended their results also to the  $p$ -adic case.

In the first part of this paper (see Theorem 2 in [5]), we derived as a consequence of our main result, effective lower bounds for linear forms with algebraic coefficients under condition (2) concerning the coefficients. Our result generalized Theorem 2 of Györy and Papp [17] to the inhomogeneous case. Our Theorem 1 makes possible to extend the special case  $L=\mathbf{Q}$  of Theorem 2 of [5] to the case when the coefficients of the linear form satisfy (3) instead of (2).

Let  $K$  be an algebraic number field of degree  $n$  and let  $r, R_K, R_K^*$  be as in Theorem 1. Denote by  $s$  and  $2t$  the number of real and complex conjugate fields of  $K$  over  $\mathbf{Q}$ . Further, denote by  $\Omega$  the set of all archimedean valuations  $|\cdot|_v$  of  $K$  where  $v=1, 2, \dots, s+t$ . For  $\beta \in K$  put  $\|\beta\|_v = |\beta|_v^{n_v}$  with  $n_v = [K_v:\mathbf{Q}_v]$ . Let  $\Gamma$  be a non-empty subset of  $\Omega$  and denote by  $s'$  and  $t'$  the number of real and complex valuations of  $\Gamma$ .

Let  $\alpha_0=1, \alpha_1, \dots, \alpha_k$  be elements in  $K$  with heights at most  $H$  ( $H \geq e$ ) such that  $K=\mathbf{Q}(\alpha_1, \dots, \alpha_k)$ . Suppose that  $\alpha_i$  is of degree  $\geq 3$  over  $\mathbf{Q}(\alpha_0, \dots, \alpha_{i-1})$  for  $i=1, \dots, k$ . Finally, let  $C$  ( $\geq e$ ) and  $0 < \zeta < 1$  be given constants.

**THEOREM 2.** Let  $x_0, x_1, \dots, x_k \in \mathbf{Z}$  and  $\lambda \in \mathbf{Z}_K$  with  $|\overline{\lambda}| < C(\max_{0 \leq i \leq k} |x_i|)^{1-\zeta}$ . If  $x_0 + x_1\alpha_1 + \dots + x_k\alpha_k + \lambda \neq 0$  then we have

$$(7) \quad \prod_{v \in \Gamma} \|x_0 + x_1\alpha_1 + \dots + x_k\alpha_k + \lambda\|_v > q_1 X^{-n+s'+2t'+\tau_1}, \quad X = \max_{0 \leq i \leq k} |x_i|$$

with

$$\varrho_1 = (H^{2n^2}C)^{-1}(4HkC)^{-n+s'+2t'}, \quad \tau_1 = \frac{\zeta^2}{9^{t-1}c_1},$$

(where  $c_1$  denotes the same expression as in Theorem 1).

Our Theorem 2 includes the case  $L=Q$  of Theorem 2 of [5]. Further, in the special case  $\lambda=0$  it implies the case  $L=Q$  of Theorem 2 of Győry and Papp [17].

**COROLLARY 2.1.** Let  $\alpha_1, \dots, \alpha_k$  be as in Theorem 2. Let  $x_1, \dots, x_k \in \mathbf{Z}$  and  $\lambda \in \mathbf{Z}_K$  with  $|\lambda| < C(\max_{1 \leq i \leq k} |x_i|)^{1-\zeta}$ . If  $\|x_1\alpha_1 + \dots + x_k\alpha_k + \lambda\| > 0$  then we have

$$(8) \quad \|x_1\alpha_1 + \dots + x_k\alpha_k + \lambda\| > \varrho_2 X^{-(n-\sigma-\tau_2)/\sigma}, \quad X = \max_{1 \leq i \leq k} |x_i|$$

where

$$\varrho_2 = (H^{2n^2}C)^{-1/\sigma}(4HkC)^{(-2n+2\sigma)/\sigma}, \quad \tau_2 = \tau_1$$

and  $\sigma=1$  or 2 according as  $K$  is real or not.

The statement of this corollary implies Corollary 2.1 of [5] and Corollary of Theorem 2 of Győry and Papp [17].

**COROLLARY 2.2.** Let  $\vartheta$  be an algebraic number of degree  $n \geq 3$  with height  $\leq H$ . Let  $K=Q(\vartheta)$  with parameters as above. Let  $p, q \in \mathbf{Z}$ ,  $q \neq 0$ , and let  $\lambda \in \mathbf{Z}_K$  with  $|\lambda| < C|q|^{1-\zeta}$ . If  $q\vartheta - p - \lambda \neq 0$  then we have

$$(9) \quad \left| \vartheta - \frac{p+\lambda}{q} \right| > \varrho_3 |q|^{-(n-\tau_3)/\sigma}$$

where we get  $\varrho_3$  and  $\tau_3$  from  $\varrho_2$  and  $\tau_2$ , respectively by taking  $k=1$ , and  $\sigma=1$  or 2 according as  $K$  is real or not.

In the special case  $\lambda=0$  Corollary 2.2 provides a slightly weaker form of a result of Győry and Papp [17] who gave an explicit version of Feldman's theorem [4].

### 3. Proofs

The proof of our Theorem 1 is based on the following Lemma:

**LEMMA.** Let  $K, \alpha_1=1, \alpha_2, \dots, \alpha_k, m$  be as in Theorem 1 and let  $1 \leq l \leq k$ . Let us consider the equation

$$(10) \quad N_{K|Q}(\alpha_1 x_1 + \dots + \alpha_l x_l + \lambda_l) = m$$

where the variables are  $x_1, \dots, x_l \in \mathbf{Z}$  and  $\lambda_l \in \mathbf{Z}_K$  with  $|\lambda_l| < C_l |x_l|^{1-\zeta_l}$  ( $C_l \geq e$  and  $0 < \zeta_l < 1$  are given constants). Then for all solutions of equation (10) we have

$$(11) \quad \max_{1 \leq i \leq l} |x_i| < \exp \left[ \frac{4n9^{l-2}}{\zeta_l} \left( \frac{D_0}{\zeta_l} + \log C_l \right) \right]$$

\*  $\| \cdot \|$  denotes the distance from the nearest integer.

with

$$D_0 = (25(r+3)n)^{17(r+3)} (R_K^* \log R_K^*)^2 (2n^2 \log H + \log |m|).$$

PROOF OF THE LEMMA. In our proof we combine the arguments of Sprindžuk [32] and Györy and Papp [17].

In the case  $l=1$  equation (10) can be written in the form

$$(12) \quad \prod_{i=1}^n (x_1 + \lambda_1^{(i)}) = m$$

where  $\lambda_1^{(i)}$  ( $i=1, \dots, n$ ) denote the conjugates of  $\lambda_1$ . Since  $|\lambda_1| < C_1 |x_1|^{1-\zeta_1}$ , (12) implies

$$|x_1|^n \leq 2^n |x_1|^{n-\zeta_1} C_1^n + |m|$$

that is

$$(13) \quad |x_1| \leq (2^n C_1^n + |m|)^{1/\zeta_1}.$$

Now let us consider the case  $l \geq 2$ . For simplicity we shall omit the index  $l$  of  $\lambda_l$ ,  $C_l$ ,  $\zeta_l$  in the first part of the proof of our Lemma.

1. First let us consider the case when

$$(14) \quad \log |x_l| > \frac{2}{\zeta} \log C + \frac{1}{\zeta} (2rn)^{6r+5} R_K^* (2n^2 \log H + \log |m|).$$

Put  $L = \mathbf{Q}(\alpha_1, \dots, \alpha_{l-1})$ ,  $[K:L] = n_L (\geq 3)$  and let  $x = \alpha_1 x_1 + \dots + \alpha_{l-1} x_{l-1}$ ,  $\vartheta = \alpha_l$  and  $y = x_l$ . With this notation equation (10) implies

$$(15) \quad N_{K|L}(x + \vartheta y + \lambda) = \mu_1$$

where  $\mu_1 \in L$  with  $N_{L|\mathbf{Q}}(\mu_1) = m$ . By Lemma 3 of Györy [7] there exists a unit  $\varepsilon$  in  $L$  such that for  $\mu_2 = \mu_1 \varepsilon^{n_L}$  we have

$$(16) \quad |\mu_2| \leq |m| \exp \left( c_2^2 \frac{r! r n^2}{2} R_K \right) = c_3$$

with  $c_2 = \left( \frac{6rn^2}{\log n} \right)^r$ . (In (16) we used also the fact that the regulator  $R_L$  of  $L$  is less than  $r! n c_2 R_K$  — see e.g. Sprindžuk [33].) From (15) we get

$$(17) \quad N_{K|L}(\varepsilon x + \vartheta \varepsilon y + \varepsilon \lambda) = \mu_2.$$

Let us consider an isomorphism  $K \rightarrow K'$  into the complex numbers and denote by  $L'$ ,  $\varepsilon'$ ,  $x'$ ,  $\vartheta'$ ,  $\lambda'$ ,  $\mu'_2$  the conjugates of  $L$ ,  $\varepsilon$ ,  $x$ ,  $\vartheta$ ,  $\lambda$ ,  $\mu_2$ , respectively under this isomorphism. Let us choose this isomorphism so that  $|\varepsilon'| \geq 1$ . By (17)  $x'$ ,  $y$  and  $\lambda'$  satisfy the equation

$$(18) \quad N_{K'|L'}(\varepsilon' x' + \vartheta' \varepsilon' y + \varepsilon' \lambda') = \mu'_2.$$

Denote by  $\vartheta'_i$ ,  $\lambda'_i$  and  $\beta'_i = x' + \vartheta'_i y + \lambda'_i$  ( $i=1, \dots, n_L$ ) the conjugates of  $\vartheta'$ ,  $\lambda'$  and  $\beta' = x' + \vartheta' y + \lambda'$ , respectively over  $L'$ . We may choose the indices so that  $|\beta'_1| \leq \dots \leq |\beta'_{n_L}|$ . Since

$$|\beta'_1 - \beta'_j| \leq |\beta'_1| + |\beta'_j| \leq 2|\beta'_j| \quad (2 \leq j \leq n_L)$$

we have

$$(19) \quad |\beta'_j| \cong \frac{1}{2} |\beta'_1 - \beta'_j| = \frac{1}{2} |\vartheta'_1 y - \vartheta'_j y + \lambda'_1 - \lambda'_j|.$$

Denote by  $a > 0$  the leading coefficient of the minimal defining polynomial of  $\vartheta'$  over  $\mathbf{Z}$ . Then  $a(\vartheta'_1 - \vartheta'_j)$  is a non-zero algebraic integer for any  $j$  with  $2 \leq j \leq n_L$  and thus

$$(20) \quad a |\vartheta'_1 - \vartheta'_j| \leq 2 \overline{|\vartheta' a|} \leq 4H$$

whence we have

$$(21) \quad |\vartheta'_1 - \vartheta'_j| \leq (4H)^{-n(n-1)} = c_4.$$

Using the above inequality we get

$$c_4 |y| \leq |\vartheta'_1 y - \vartheta'_j y| \leq |\vartheta'_1 y - \vartheta'_j y + \lambda'_1 - \lambda'_j| + |\lambda'_1 - \lambda'_j|.$$

Further, applying (19) and (14) we have

$$(22) \quad |\beta'_j| \cong \frac{1}{2} |\vartheta'_1 y - \vartheta'_j y + \lambda'_1 - \lambda'_j| \cong \frac{c_4}{2} |y| - C |y|^{1-\zeta} \cong \frac{c_4}{4} |y| \quad (2 \leq j \leq n_L).$$

Equation (18) implies

$$\prod_{j=1}^{n_L} (\varepsilon' \beta'_j) = \mu'_2$$

and thus applying (16), from (22) we get

$$|\varepsilon' \beta'_1| = \frac{|\mu'_2|}{\prod_{j=2}^{n_L} |\varepsilon' \beta'_j|} \leq c_3 \left( \frac{c_4}{4} |\varepsilon' y| \right)^{1-n_L} = c_5 |\varepsilon' y|^{1-n_L}.$$

Since  $|\varepsilon'| \geq 1$  the above inequality implies

$$(23) \quad |\beta'_1| \leq c_5 |y|^{1-n_L}.$$

By (23) we get for any  $j$ ,  $1 \leq j \leq n_L$ ,

$$(24) \quad |\beta'_j| \leq |\beta'_j - \beta'_1| + |\beta'_1| \leq |(\vartheta'_j - \vartheta'_1)y + \lambda'_j - \lambda'_1| + c_5 |y| \leq c_6 |y|.$$

By Lemma 2 of Györy [7] there exist fundamental units  $\eta_1, \dots, \eta_r$  in  $K$  such that

$$(25) \quad \prod_{j=1}^r \max(\log |\overline{\eta_j}|, 1) < c_2 R_K$$

and the absolute values of the elements of the inverse matrix of  $(e_i \log |\eta_j^{(i)}|)_{1 \leq i, j \leq r}$  do not exceed  $c_7$  where  $e_i = 1$  or  $2$  according as the conjugate field  $K^{(i)}$  of  $K$  is real or complex and  $c_7 = \frac{6r!n^2}{\log n}$ .

Put  $\beta = x + \vartheta y + \lambda$ . Since  $|N_{K|Q}(\beta)| = m$ , by Lemma 3 of Györy [7] there exist rational integers  $b_1, \dots, b_r$  such that for  $\gamma = \beta \eta_1^{b_1} \dots \eta_r^{b_r}$  we have

$$(26) \quad \left| \log \left| |m|^{-1/n} \gamma^{(j)} \right| \right| \leq \frac{c_2 r}{2} R_K \quad (1 \leq j \leq n).$$

Denote by  $\eta'_{i,j}$  and  $\gamma'_j$  the conjugates of  $\eta_i$  and  $\gamma$  corresponding to  $\beta'_j$ . By (14), (24) we have

$$\begin{aligned} |b_1 \log |\eta'_{1,j}| + \dots + b_r \log |\eta'_{r,j}| &= \left| \log \left| \frac{\gamma'_j}{\beta'_j} \right| \right| = \\ &= \left| \log \left| |m|^{-1/n} \gamma'_j \right| + \frac{1}{n} \log |m| - \log |\beta'_j| \right| \leq \\ &\leq \frac{c_2 r}{2} R_K + \frac{1}{n} \log |m| + \log c_6 + \log |y| \leq \frac{3}{2} \log |y| \end{aligned}$$

and from this inequality we get

$$(27) \quad \max_{1 \leq i \leq r} |b_i| \leq 3rc_7 \log |y|.$$

Since  $n_L \equiv 3$  we may suppose that  $\vartheta'_1, \vartheta'_2, \vartheta'_3$  are pairwise distinct. Let us consider now the following identity:

$$(\vartheta'_2 - \vartheta'_3)(\beta'_1 - \lambda'_1) + (\vartheta'_3 - \vartheta'_1)(\beta'_2 - \lambda'_2) + (\vartheta'_1 - \vartheta'_2)(\beta'_3 - \lambda'_3) = 0.$$

Applying our estimates (21), (22), (23) and using (14) from this identity we obtain

$$\begin{aligned} \left| 1 + \frac{(\vartheta'_1 - \vartheta'_2)\beta'_3}{(\vartheta'_3 - \vartheta'_1)\beta'_2} \right| &\leq \left| \frac{(\vartheta'_2 - \vartheta'_3)\beta'_1}{(\vartheta'_3 - \vartheta'_1)\beta'_2} \right| + \left| \frac{(\vartheta'_2 - \vartheta'_3)\lambda'_1 + (\vartheta'_3 - \vartheta'_1)\lambda'_2 + (\vartheta'_1 - \vartheta'_2)\lambda'_3}{(\vartheta'_3 - \vartheta'_1)\beta'_2} \right| \leq \\ &\leq \frac{4H}{c_4} \frac{c_5 |y|^{1-n_L} + 3C |y|^{1-\zeta}}{\frac{c_4}{4} |y|} \leq \frac{1}{2} |y|^{-\zeta/2}. \end{aligned}$$

Let

$$\varrho_i = \begin{cases} \frac{\eta'_{i,2}}{\eta'_{i,3}} & \text{for } i = 1, \dots, r \\ \frac{(\vartheta'_2 - \vartheta'_1)\gamma'_3}{(\vartheta'_3 - \vartheta'_1)\gamma'_2} & \text{for } i = r+1. \end{cases}$$

With this notation the above inequality can be written in the form

$$(28) \quad 0 < |\varrho_1^{b_1} \dots \varrho_r^{b_r} \varrho_{r+1} - 1| < \frac{1}{2} |y|^{-\zeta/2}.$$

Put  $\varrho_0 = -1$ . Then (28) implies

$$(29) \quad 0 < |b_0 \log \varrho_0 + b_1 \log \varrho_1 + \dots + b_r \log \varrho_r - \log \varrho_{r+1}^{-1}| < e^{-(\zeta/2) \log |y|} = e^{-\delta B}$$

where  $b_0 \in \mathbf{Z}$  with  $|b_0| \leq |b_1| + \dots + |b_r|$ ,  $\log$  denotes the principal value of the logarithm, by (27) we have

$$\max_{0 \leq i \leq r} |b_i| \leq 3r^2 c_7 \log |y| = B$$

and  $\delta = \zeta (6r^2 c_7)^{-1}$ . Let  $A_i = \max(H(\varrho_i), e^e)$  for  $i=0, 1, \dots, r$ . For any  $i$ ,  $1 \leq i \leq r$  we have

$$H(\varrho_i) \leq (2 \overline{|\eta'_{i,2}/\eta'_{i,3}|})^{n(n-1)} \leq (2 \overline{|\eta_i|})^{n(n-1)}$$

that is

$$\log A_i \leq 2n^2(n-1) \max(\log \overline{|\eta_i|}, 1) \quad (1 \leq i \leq r).$$

Applying (25), from this inequality we get

$$(30) \quad \Omega' = \log A_0 \log A_1 \dots \log A_r \leq 4(2n^2(n-1))^r c_2 R_K = c_8.$$

Denote by  $a_i$  the leading coefficient of the minimal defining polynomial of  $\alpha_i$  over  $\mathbf{Z}$ , where  $\alpha_i$  ( $1 \leq i \leq l$ ) denotes the coefficients in equation (10). Further, put  $a^* = a_2 \dots a_l$ . Then  $a^* \gamma$  is an algebraic integer and from (20), (26) we obtain

$$\begin{aligned} H(\varrho_{r+1}) &\leq (\overline{|a(\vartheta'_2 - \vartheta'_1) a^* \gamma'_3|} + \overline{|a(\vartheta'_3 - \vartheta'_1) a^* \gamma'_2|})^{n(n-1)(n-2)} \leq \\ &\leq [8H^k |m|^{1/m} \exp(2c_2 r R_K)]^{n^2(n-2)} = A. \end{aligned}$$

For the above  $A$  we have  $A_i \leq A$ ,  $0 \leq i \leq r$ . Put  $c_9 = (25(r+3)n)^{10(r+3)}$  and  $T = c_9 \Omega' \log \Omega'$ . Obviously,  $\delta \leq c_9^{-1/2} T$ . Now we apply Theorem 3 of van der Poorten and Loxton [26] (see also [27]) to (29) and we obtain

$$B < \delta^{-1} T \log(\delta^{-1} T) \log A.$$

Applying our estimate (30), from the above inequality we get an upper bound for  $|y| = |x_l|$ . Combining this bound with the bound we get for  $|x_l|$  in the case when (14) does not hold, we obtain

$$|y| < \exp \left[ \frac{D_0}{\zeta^2} + \frac{2}{\zeta} \log C \right],$$

that is, with the notation of our Lemma

$$(31) \quad |x_l| < \exp \left[ \frac{D_0}{\zeta_l^2} + \frac{2}{\zeta_l} \log C_l \right] = D_l.$$

If  $l=2$ , we may continue immediately at part III.

II. If  $l > 2$ , put  $L = \mathbf{Q}(\alpha_1, \dots, \alpha_{l-2})$ ,  $x = \alpha_1 x_1 + \dots + \alpha_{l-2} x_{l-2}$ ,  $\vartheta = \alpha_{l-1}$ ,  $y = x_{l-1}$  and  $\lambda_{l-1} = \alpha_l x_l + \lambda_l$ . Then  $x, y, \lambda_{l-1}$  satisfy

$$(32) \quad N_{K|L}(x + \vartheta y + \lambda_{l-1}) = \mu_1$$

where  $\mu_1 \in L$  with  $|N_{L|\mathbf{Q}}(\mu_1)| = |m|$ . Further,  $|\overline{\lambda_{l-1}}| < (2Hk + C_l) D_l = C_{l-1}$ . Applying the arguments of part I to equation (32) with  $l-1$  instead of  $l$ , with  $C_{l-1}$ ,  $\zeta_{l-1} = 1/2$  and with the above  $\lambda_{l-1}$  we obtain

$$\begin{aligned} |x_{l-1}| &< \exp [4D_0 + 4 \log C_{l-1}] = \exp [4D_0 + 4 \log (2Hk + C_l) + 4 \log D_l] < \\ &< \exp (9 \log D_l) = D_{l-1}. \end{aligned}$$



Continuing this procedure we obtain

$$(33) \quad \max_{2 \leq i \leq l} |x_i| < \exp [9^{l-2} \log D_l] = D_2.$$

III. Finally, put  $\lambda_1 = \alpha_2 x_2 + \dots + \alpha_l x_l + \lambda_l$ . Then  $|\lambda_1| < (2Hk + C_l)D_2$  and  $x_1, \lambda_1$  satisfy equation

$$(34) \quad \prod_{i=1}^n (x_1 + \lambda_1^{(i)}) = m$$

where  $\lambda_1^{(i)}$  ( $1 \leq i \leq n$ ) denote the conjugates of  $\lambda_1$  over  $\mathbf{Q}$ . Assuming  $x_1 \neq 0$ , from (34) we obtain

$$|x_1| < 2^n [(2Hk + C_l)D_2]^n + |m|$$

that is

$$|x_1| < \exp (2n \log D_2).$$

Combining this with (31) and (33), in the case  $l > 1$  we get

$$\max_{1 \leq i \leq l} |x_i| < \exp [2n 9^{l-2} \log D_l] \leq \exp \left[ \frac{2n 9^{l-2}}{\zeta_l} \left( \frac{D_0}{\zeta_l} + 2 \log C_l \right) \right].$$

This proves (11) in the case  $l > 1$ , and as we have seen in (13), in the case  $l = 1$  one can get a much better bound for  $|x_1|$ .

PROOF OF THEOREM 1. We can easily prove Theorem 1 by applying our Lemma.

Let us consider an arbitrary solution  $x_1, \dots, x_k, \lambda$  of equation (5) with  $|\lambda| < CX^{1-\zeta}$ , where  $X = \max_{1 \leq i \leq k} |x_i|$ . Put  $\mathcal{L}(\mathbf{x}, \lambda) = \alpha_1 x_1 + \dots + \alpha_k x_k + \lambda$ . Denote by  $j$  the greatest index such that  $|x_j| = X$ . If  $j = k$ , then our Lemma can be directly applied with  $l = k$ ,  $\lambda_k = \lambda$ ,  $C_k = C$  and  $\zeta_k = \zeta$ . If  $j < k$  then two cases are possible.

If  $|x_{j+1}|, \dots, |x_k|$  are less than  $X^{1-(\zeta/k)}$ , put  $\lambda_j = \mathcal{L}(\mathbf{x}, \lambda) - (\alpha_1 x_1 + \dots + \alpha_j x_j)$ . Since in this case  $|\lambda_j| < c_{10} X^{1-(\zeta/k)} = c_{10} |x_j|^{1-(\zeta/k)}$  with  $c_{10} = 2kH + C$ , applying our Lemma with  $l = j$ ,  $C_j = c_{10}$ ,  $\zeta_j = \zeta/k$  and with the above  $\lambda_j$  we get the upper bound required for  $\max_{1 \leq i \leq j} |x_i| = X$ .

Let us consider now the case when there exists an index  $j' > j$  such that  $|x_{j'}| > X^{1-(\zeta/k)}$ . Denote by  $j_1$  the greatest index with this property. There are again two possible cases.

If  $j_1 = k$  or if  $j_1 < k$  and  $|x_{j_1+1}|, \dots, |x_k|$  are less than  $X^{1-(2\zeta/k)}$ , put  $\lambda_{j_1} = \mathcal{L}(\mathbf{x}, \lambda) - (\alpha_1 x_1 + \dots + \alpha_{j_1} x_{j_1})$ . Since  $|x_{j_1}| > X^{1-(\zeta/k)}$  we have  $X < |x_{j_1}|^{1/(1-(\zeta/k))}$  and thus  $|\lambda_{j_1}| < c_{10} X^{1-(2\zeta/k)} < c_{10} |x_{j_1}|^{(1-(2\zeta/k))/(1-(\zeta/k))} \leq c_{10} |x_{j_1}|^{1-(\zeta/k)}$ . Applying our Lemma with  $l = j_1$ ,  $C_{j_1} = c_{10}$ ,  $\zeta_{j_1} = \zeta/k$  and with the above  $\lambda_{j_1}$  we get the upper bound announced for  $\max_{1 \leq i \leq j_1} |x_i| = X$ .

If  $j_1 < k$  and there is an index  $j'' > j_1$  such that  $|x_{j''}| > X^{1-(2\zeta/k)}$  then denote by  $j_2$  the greatest index with this property and continue the above procedure. In each step, either we apply our Lemma with some  $l \leq k$ ,  $C_l \leq c_{10}$ ,  $\zeta_l \leq \zeta/k$ , get an upper bound for  $X$ , and stop the procedure, or we go on to the next step. Since  $j < j_1 < j_2 < \dots$ , this procedure stops in  $k$  steps at most and our theorem is proved.

PROOF OF THEOREM 2. In our proof we follow the arguments of Györy and Papp [17]. Consider arbitrary  $x_0, x_1, \dots, x_k \in \mathbf{Z}$  and  $\lambda \in \mathbf{Z}_K$  such that  $|\lambda| < CX^{1-\zeta}$

with  $X = \max_{0 \leq i \leq k} |x_i|$  and  $x_0 + x_1 \alpha_1 + \dots + x_k \alpha_k + \lambda \neq 0$ , and put

$$(35) \quad m = N_{K|Q}(x_0 + x_1 \alpha_1 + \dots + x_k \alpha_k + \lambda).$$

Applying our Theorem 1 to the above equation we get

$$X < (H^{2n^2} |m| C)^{1/\tau_1}$$

with the constant  $\tau_1$  of Theorem 2, that is

$$(36) \quad |m| > (H^{2n^2} C)^{-1} X^{\tau_1}.$$

By equation (35) we have

$$(37) \quad \prod_{v \in \Omega} \|x_0 + x_1 \alpha_1 + \dots + x_k \alpha_k + \lambda\|_v = |m|.$$

Further

$$\|x_0 + x_1 \alpha_1 + \dots + x_k \alpha_k + \lambda\|_v \leq (2H(k+1)X + CX)^{n_v} \leq (4HkCX)^{n_v}$$

for each valuation  $v \in \Omega$ . Thus, from (36) and (37) we obtain

$$\begin{aligned} \prod_{v \in \Gamma} \|x_0 + x_1 \alpha_1 + \dots + x_k \alpha_k + \lambda\|_v &\geq \frac{|m|}{(4HkCX)^{n-s'-2t'}} > \\ &> (H^{2n^2} C)^{-1} (4HkC)^{-n+s'+2t'} X^{-n+s'+2t'+\tau_1} \end{aligned}$$

which proves (7).

**PROOF OF COROLLARY 2.1.** Denote by  $-y$  the nearest integer to  $x_1 \alpha_1 + \dots + x_k \alpha_k + \lambda$ . Then we get

$$\|x_1 \alpha_1 + \dots + x_k \alpha_k + \lambda\| = |y + x_1 \alpha_1 + \dots + x_k \alpha_k + \lambda| > 0.$$

If  $\|x_1 \alpha_1 + \dots + x_k \alpha_k + \lambda\| \geq 1$  then (8) obviously holds. If  $\|x_1 \alpha_1 + \dots + x_k \alpha_k + \lambda\| < 1$  then  $|y| < 1 + 2HkX + CX < 4HkCX$ . Applying our Theorem 2 to the inhomogeneous linear form  $y + x_1 \alpha_1 + \dots + x_k \alpha_k + \lambda$ , we get (8).

**PROOF OF COROLLARY 2.2.** Applying Corollary 2.1 in the special case  $k=1$  we obtain

$$\|q\vartheta - \lambda\| > \varrho_3 |q|^{-(n-\sigma-\tau_3)/\sigma},$$

that is, for any  $p \in \mathbb{Z}$  we have

$$|-p + q\vartheta - \lambda| > \varrho_3 |q|^{-(n-\sigma-\tau_3)/\sigma}.$$

Dividing this inequality by  $|q|$  we get (9).

**ACKNOWLEDGEMENTS.** I am thankful to Professor Kálmán Györy for his help in the preparation of this paper.

#### REFERENCES

- [1] BAKER, A., Contributions to the theory of Diophantine equations, *Philos. Trans. Roy. Soc. London Ser. A* **263** (1968), 173—191. *MR* **37** #4005.
- [2] BAKER, A., *Transcendental number theory*, Univ. Press, Cambridge, 1979. *MR* **54** #10163.
- [3] COATES, J., An effective  $p$ -adic analogue of a theorem of Thue, *Acta Arith.* **15** (1969), 279—305. *MR* **39** #4095.

- [4] FELDMAN, N. I., An effective power sharpening of a theorem of Liouville (in Russian), *Izv. Akad. Nauk SSSR* **35** (1971), 973—990. *MR* **44**#6609.
- [5] GAÁL, I., Norm form equations with several dominating variables and explicit lower bounds for inhomogeneous linear forms with algebraic coefficients, *Studia Sci. Math. Hungar.* to appear.
- [6] GYÖRY, K., On the greatest prime factors of decomposable forms at integer points, *Ann. Acad. Sci. Fenn. Ser. A I* **4** (1978/1979), 341—355. *MR* **81g**: 10038.
- [7] GYÖRY, K., On the solutions of linear diophantine equations in algebraic integers of bounded norm, *Ann. Univ. Budapest Eötvös Sect. Math.* **22/23** (1979/1980), 225—233. *MR* **82c**: 10016.
- [8] GYÖRY, K., Explicit lower bounds for linear forms with algebraic coefficients, *Arch. Math. (Basel)* **35** (1980), 438—446. *MR* **82c**: 10043.
- [9] GYÖRY, K., Explicit upper bounds for the solutions of some diophantine equations, *Ann. Acad. Sci. Fenn. Ser. A I. Math.* **5** (1980), 3—12. *MR* **82e**: 10028.
- [10] GYÖRY, K., Sur certaines généralisations de l'équation de Thue—Mahler, *Enseign. Math.* **26** (1980), 247—255. *MR* **82e**: 10032.
- [11] GYÖRY, K., *Résultats effectifs sur la représentation des entiers par des formes décomposables*, Queen's Papers in Pure and Applied Math., No. 56, Kingston (Canada), 1980. *MR* **83c**: 10021.
- [12] GYÖRY, K., On the representation of integers by decomposable forms in several variables, *Publ. Math. Debrecen* **28** (1981), 89—98. *MR* **83b**: 10017.
- [13] GYÖRY, K., On  $S$ -integral solutions of norm form, discriminant form and index form equations, *Studia Sci. Math. Hungar.* **16** (1981), 149—161.
- [14] GYÖRY, K., Bounds for the solutions of norm form, discriminant form and index form equations in finitely generated integral domains, *Acta Math. Acad. Sci. Hungar.* **42** (1983), 45—80.
- [15] GYÖRY, K., On norm form, discriminant form and index form equations, *Colloquia Math. Soc. János Bolyai* **34. Topics in Classical Number Theory**, Budapest (Hungary), 1981, 617—676 (1984).
- [16] GYÖRY, K. and PAPP, Z. Z., Effective estimates for the integer solutions of norm form and discriminant form equations, *Publ. Math. Debrecen* **25** (1978), 311—325. *MR* **80b**: 10026.
- [17] GYÖRY, K. and PAPP, Z. Z., Norm form equations and explicit lower bounds for linear forms with algebraic coefficients, *Studies in Pure Mathematics (To the memory of Paul Turán)*, Akadémiai Kiadó, Budapest, 1983, 245—257.
- [18] KOTOV, S. V., The Thue—Mahler equation in relative fields (in Russian), *Acta Arith.* **27** (1975), 293—315. *MR* **51**#12722.
- [19] KOTOV, S. V., On diophantine equations of norm form type I. (in Russian), *Inst. Mat. Akad. Nauk BSSR*, Preprint No. 9, Minsk, 1980.
- [20] KOTOV, S. V., On diophantine equations of norm form type II. (in Russian), *Inst. Mat. Akad. Nauk BSSR*, Preprint No. 10, Minsk, 1980.
- [21] KOTOV, S. V., Effective bounds for linear forms with algebraic coefficients in archimedean and  $p$ -adic metrics (in Russian), *Inst. Mat. Akad. Nauk BSSR*, Preprint No. 24 Minsk 1981.
- [22] KOTOV, S. V., Effective bound for the values of the solutions of a class of diophantine equations of norm form type (in Russian), *Mat. Zametki* **33** (1983), 801—806.
- [23] KOTOV, S. V. and SPRINDŽUK, V. G., The Thue—Mahler equation in relative fields and approximation of algebraic numbers by algebraic numbers (in Russian), *Izv. Akad. Nauk SSSR* **41** (1977), 723—751. *MR* **58**#5539.
- [24] LIOUVILLE, J., Sur des classes très étendues de quantités dont la valeur n'est ni algébrique, ni même réductible à des irrationnelles algébriques, *C. R. Acad. Sci. Paris* **18** (1844), 883—885 and 910—911.
- [25] MAHLER, K., Zur Approximation algebraischer Zahlen I: Über den grössten Primteiler binärer Formen, *Math. Ann.*, **107** (1933), 691—730.
- [26] VAN DER POORTEN, A. J. and LOXTON, J. H., Multiplicative relations in number fields, *Bull. Austral. Math. Soc.* **16** (1977), 83—98. *MR* **58**#10776a.
- [27] VAN DER POORTEN, A. J. and LOXTON, J. H., Computing the effectively computable bound in Baker's inequality for linear forms in logarithms, and, Multiplicative relations in number fields: Corrigendum and addendum, *Bull. Austral. Math. Soc.* **17** (1977), 151—155. *MR* **58**#10776b.

- [28] SHOREY, T. N., VAN DER POORTEN, A. J., TIJDEMAN, R. and SCHINZEL, A., Applications of the Gelfond—Baker method to Diophantine equations, *Transcendence Theory: Advances and Applications* (ed. by A. Baker and D. W. Masser), Academic Press, London—New York—San Francisco, 1977, 59—77. *MR* 57#12383.
- [29] SPRINDŽUK, V. G., A new application of  $p$ -adic analysis to representation of numbers by binary forms (in Russian), *Izv. Akad. Nauk SSSR* 34 (1970), 1038—1063. *MR* 42#5910.
- [30] SPRINDŽUK, V. G., Rational approximations to algebraic numbers (in Russian), *Izv. Akad. Nauk SSSR* 35 (1971), 991—1007. *MR* 45#1846.
- [31] SPRINDŽUK, V. G., Estimation of the solutions of the Thue equation (in Russian), *Izv. Akad. Nauk SSSR* 36 (1972), 712—741. *MR* 47#1741.
- [32] SPRINDŽUK, V. G., Representation of numbers by the norm forms with two dominating variables, *J. Number Theory* 6 (1974), 481—486. *MR* 50#7045.
- [33] SPRINDŽUK, V. G., *Classical diophantine equations in two variables* (in Russian), Izd. Nauka, Moscow, 1982.
- [34] STARK, H. M., Effective estimates of solutions of some Diophantine equations, *Acta Arith.* 24 (1973), 251—259. *MR* 49#4931.
- [35] THUE, A., Über Annäherungswerte algebraischer Zahlen, *J. Reine Angew. Math.* 135 (1909), 284—305.

( Received October 25, 1983 )

KOSSUTH LAJOS TUDOMÁNYEGYETEM  
MATEMATIKAI INTÉZETE  
P.O. BOX 12  
H-4010 DEBRECEN  
HUNGARY

# ON $2k$ -DIMENSIONAL DENSITY ESTIMATES

ALBERTO PERELLI and SAVERIO SALERNO

## 1. Introduction

Density theorems for the zeros of zeta and  $L$  functions find applications to several problems concerning the distribution of prime numbers. They consist in estimates for the density functions  $N(\sigma, T, x)$  and are generally used as quantitative substitutes of the conjectural assertion that in fact  $L(s, x) \neq 0$  for  $\sigma > 1/2$ .

Recently Heath-Brown [1], [2] introduced a new type of density function and applied his estimates to the problem of the second moment of the differences between consecutive primes. Precisely, he introduced the quantity

$$N^*(\sigma, T) = |\{\varrho_j = \beta_j + i\gamma_j: \zeta(\varrho_j) = 0, \beta_j \leq \sigma, |\gamma_j| \leq T, |\gamma_1 + \gamma_2 - \gamma_3 - \gamma_4| \leq 1\}|$$

and his techniques were based on an ingenious variant of the classical zero-detection method.

Our aim in this paper is to give some estimates for the  $2k$ -dimensional density function

$$(1.1) \quad N_{2k}(\sigma, T) = |\{\varrho_j: \zeta(\varrho_j) = 0, \beta_j \leq \sigma, |\gamma_j| \leq T, j = 1, 2, \dots, 2k, \\ |\gamma_1 + \dots + \gamma_k - \gamma_{k+1} - \dots - \gamma_{2k}| \leq 1\}|.$$

For  $k=1, 2$  we have

$$N_2(\sigma, T) \ll N(\sigma, T) \log T, \quad N_4(\sigma, T) = N^*(\sigma, T)$$

and our method for general  $k$  is based on Heath-Brown's one in [1] and [2].

The trivial estimate for  $N_{2k}(\sigma, T)$  is

$$(1.2) \quad N_{2k}(\sigma, T) \ll N(\sigma, T)^{2k-1} \log T;$$

we obtain non-trivial estimates only in a neighbourhood of  $\sigma=3/4$ , which is often the critical point in the applications, and for  $k > k_0$ .

Our result is the following

---

1980 *Mathematics Subject Classification*. Primary 11M26; Secondary 11N05.  
*Key words and phrases*. Density theorems.

THEOREM 1. For every  $\varepsilon > 0$  we have

$$(1.3) \quad N_{2k}(\sigma, T) \ll T^{\frac{3(1-\sigma)}{7\sigma-4} \left(2k - \frac{\sigma}{2(1-\sigma)(2\sigma-1)}\right) + \frac{1}{2(2\sigma-1)} + O\left(\frac{1}{k}\right) + \varepsilon}$$

$$\text{for } \frac{3}{4} \leq \sigma \leq \frac{6 + \sqrt{22}}{14} + O\left(\frac{1}{k}\right);$$

$$(1.4) \quad N_{2k}(\sigma, T) \ll T^{\frac{3(1-\sigma)}{2-\sigma} \left(2k - \frac{\sigma}{2(1-\sigma)(2\sigma-1)}\right) + \frac{1}{2(2\sigma-1)} + O\left(\frac{1}{k}\right) + \varepsilon}$$

$$\text{for } \frac{2}{3} + O\left(\frac{1}{k}\right) \leq \sigma \leq \frac{3}{4}.$$

We will use the density estimates of Jutila [7] and Ivić [6]; it is of course possible to use any other density estimate, eventually modifying the range for  $\sigma$ . In order to avoid complicated details we do not try to find the optimal choice of the parameters in the proof of Theorem 1.

$2k$ -dimensional density theorems may be used in order to estimate  $2k$ th-moments of primes in short intervals. However, we will show in the Appendix how such results may be obtained by means of a simple direct argument.

## 2. Proof of Theorem 1

### a) THE ZERO DETECTION METHOD

Let us briefly recall the zero-detection method, as modified by Heath-Brown [1] and [2], adapted for our general case.

Let  $X, Y$  be such that  $3 \leq X, Y \leq T^2$  and let  $M_X(s) = \sum_{n \leq X} \mu(n)n^{-s}$ . Let  $\mathcal{N}(\sigma, T)$  be the set of zeros counted by  $\mathcal{N}(\sigma, T)$  and  $\mathcal{S}_0^{(0)} = \{\varrho \in \mathcal{N}(\sigma, T) : \text{either } |\gamma| \leq (\log T)^2 \text{ or}$

$$\int_{-\infty}^{+\infty} \left| \zeta\left(\frac{1}{2} + it\right) M_X\left(\frac{1}{2} + it\right) \right| e^{-|\gamma-t|} dt \gg Y^{\beta - (1/2)}\},$$

$$\mathcal{S} = \mathcal{N}(\sigma, T) \setminus \mathcal{S}_0^{(0)}.$$

The zeros of  $\mathcal{S}$  may be divided into  $O(\log T)$  subsets  $\mathcal{S}^{(m)}$ , for which

$$\left| \sum_{n \leq X} a_n n^{-\varrho} \right| \gg (\log T)^{-1}, \quad \varrho \in \mathcal{S}^{(m)},$$

$$(2.1) \quad X/2 \leq N = 2^m \leq Y(\log T)^2,$$

where  $a_n$  are suitable coefficients satisfying  $|a_n| \leq d(n)$ .  $N_{2k}(\sigma, T)$  is then connected with the above quantities in the following way. Let

$$\mathcal{S}_j^{(m)} = \{\varrho \in \mathcal{S}^{(m)} : [\gamma] \equiv j \pmod{2k}\}, \quad 1 \leq j \leq 2k,$$

$$\mathcal{S}(x) = \sum_{\varrho \in \mathcal{N}(\sigma, T)} e(\gamma x), \quad \mathcal{S}_j^{(m)}(x) = \sum_{\varrho \in \mathcal{S}_j^{(m)}} e(\gamma x),$$

$$w(x) = \left( \frac{\sin(2\pi x)}{2\pi x} \right)^2$$



where  $e(x) = e^{2\pi i x}$  since

$$\int_{-\infty}^{+\infty} e(yx) w(x) dx = \begin{cases} 2-|y| & \text{if } |y| \leq 2 \\ 0 & \text{if } |y| > 2, \end{cases}$$

we have

$$(2.2) \quad N_{2k}(\sigma, T) \leq \int_{-\infty}^{+\infty} |\mathcal{S}(x)|^{2k} w(x) dx \ll (\log T)^{2k-1} \sum_{m,j} \int_{-\infty}^{+\infty} |\mathcal{S}_j^{(m)}(x)|^{2k} w(x) dx \ll$$

$$\ll (\log T)^{2k} \max_{m,j} \left( \sum_{\substack{q_1, \dots, q_{2k} \in \mathcal{S}_j^{(m)} \\ |\gamma_1 + \dots + \gamma_k - \gamma_{k+1} - \dots - \gamma_{2k}| \leq 2}} 1 \right).$$

The contribution of  $\mathcal{S}_0^{(0)}$  is estimated trivially: for each of  $|\mathcal{S}_0^{(0)}|^{2k-1}$  choices of  $q_1, \dots, q_{2k-1}$  there are  $O(\log T)$  choices for  $q_{2k}$ , hence the contribution of  $\mathcal{S}_0^{(0)}$  to (2.2) is  $O((\log T)^{2k+1} |\mathcal{S}_0^{(0)}|^{2k-1})$ . In order to treat the other  $\mathcal{S}_j^{(m)}$  we note that  $|\gamma_1 + \dots + \gamma_k - \gamma_{k+1} - \dots - \gamma_{2k}| \leq 2$  implies  $[\gamma_1] + \dots + [\gamma_k] - [\gamma_{k+1}] - \dots - [\gamma_{2k}] \equiv 0 \pmod{2k}$  (we have  $[\gamma_j] \equiv j \pmod{2k}$ ), hence from (2.2) we get

$$(2.3) \quad N_{2k}(\sigma, T) \ll (\log T)^{2k+1} |\mathcal{S}_0^{(0)}|^{2k-1} + (\log T)^{2k} \left( \sum_{\substack{q_1, \dots, q_{2k} \in \mathcal{S}_j^{(m)} \\ [\gamma_1] + \dots + [\gamma_k] = [\gamma_{k+1}] + \dots + [\gamma_{2k}]} 1 \right)$$

for some  $m, j$  with  $m \neq 0$ .

Let  $t_i$  run over all the values of  $[\gamma]$  for which  $q \in \mathcal{S}_j^{(m)}$ ; since the number of solutions of  $[\gamma] = t_i$ , for  $t_i$  fixed, is  $O(\log T)$ , we have

$$(2.4) \quad N_{2k}(\sigma, T) \ll (\log T)^{4k} \left( \sum_{t_1 + \dots + t_k = t_{k+1} + \dots + t_{2k}} 1 \right) + (\log T)^{2k+1} |\mathcal{S}_0^{(0)}|^{2k-1}.$$

Raising the inequality  $\left| \sum_N^{2N} a_n n^{-q} \right| \leq \frac{K}{\log T}$  to the power  $h = h(N)$  we have

$$(2.5) \quad \left| \sum_M^P b_n n^{-q} \right| \leq \left( \frac{K}{\log T} \right)^h$$

where  $M = N^h$ ,  $P = (2N)^h$ ,  $|b_n| \leq d_{2h}(n)$ . Putting  $[\gamma] = t_i = t$ ,

$$q = \sigma + it + \delta, \delta = (\beta - \sigma) + i(\gamma - [\gamma]), D(y, t) = D(y) = \sum_{M \leq n \leq y} b_n n^{-\sigma - it},$$

we get by partial summation, using (2.5):

$$(2.6) \quad 1 \ll \left( \frac{K}{\log T} \right)^{-h} \left( |D(P, t_i)| + \int_M^P |D(y, t_i)| \frac{dy}{y} \right)$$

for every  $t_i$  in (2.4).

We now let

$$R_{2k} = \sum_{t_1 + \dots + t_k = t_{k+1} + \dots + t_{2k}} 1 \quad (\text{as in (2.4)}),$$

$$m(t) = \sum_{t = t_1 + \dots + t_k - t_{k+1} - \dots - t_{2k-1}} 1, \quad n(t) = \sum_{t = t_1 + \dots + t_k} 1.$$



Then, from (2.6), we have

$$R_{2k} = \sum_{t_{2k}} m(t_{2k}) \ll \left( \frac{\log T}{K} \right)^{2h+1} \sum_{t_{2k}} m(t_{2k}) \left[ |D(P, t_{2k})|^2 + \int_M^P |D(y, t_{2k})|^2 \frac{dy}{y} \right],$$

and it is clearly sufficient to consider the first term in the sum inside square brackets. We have

$$\begin{aligned} R_{2k} &\ll \left( \frac{\log T}{K} \right)^{2h+1} \sum_{t_{2k}} \sum_{\substack{t, t' \\ t-t'=t_{2k}}} n_k(t) n_{k-1}(t') \left| \sum_M^P b_n n^{-\sigma-i(t-t')} \right|^2 \ll \\ (2.7) \quad &\ll \left( \frac{\log T}{K} \right)^{2h+1} \sum_{n_1, n_2=M}^P |b_{n_1} \bar{b}_{n_2}| (n_1 n_2)^{-\sigma} \left| \sum_t n_k(t) \left( \frac{n_2}{n_1} \right)^{it} \right| \left| \sum_{t'} n_{k-1}(t') \left( \frac{n_1}{n_2} \right)^{it'} \right| \ll \\ &\ll \left( \frac{\log T}{K} \right)^{2h+1} G (\sum_k \sum_{k-1})^{1/2} \end{aligned}$$

where

$$(2.8) \quad G = \max_{M \leq n \leq P} (d_{2k}(n))^2 n^{1-2\sigma} \ll_\varepsilon T^\varepsilon P^{1-2\sigma}$$

$$(2.9) \quad \sum_k = \sum_{n_1, n_2=M}^P (n_1 n_2)^{-1/2} \left| \sum_t n_k(t) \left( \frac{n_2}{n_1} \right)^{it} \right|^2.$$

b) THE ESTIMATE OF  $|\mathcal{S}_0^{(0)}|$ .

From the definition of  $\mathcal{S}_0^{(0)}$  we see that there are  $O((\log T)^3)$  possibilities for  $\gamma$  in the first case.

In the second case we raise the integral on the left of the inequality in the definition of  $\mathcal{S}_0^{(0)}$  to the power  $H$  and, applying Hölder's inequality, we get

$$\int_{\gamma-(\log T)^2}^{\gamma+(\log T)^2} |\zeta(1/2+it)|^H dt \gg Y^{H(\sigma-(1/2))} T^{-\varepsilon};$$

summing over  $\gamma$  we obtain

$$(2.10) \quad \int_{-T-(\log T)^2}^{T+(\log T)^2} |\zeta(1/2+it)|^H s(t) dt \gg |\mathcal{S}_0^{(0)}| Y^{H(\sigma-(1/2))} T^{-\varepsilon}$$

where

$$s(t) = |\{ \varrho \in \mathcal{S}_0^{(0)} : |t-\gamma| < \log^2 T \}| \ll (\log T)^3.$$

Using the estimates for the 4- and 12- power moment of  $\zeta(1/2+it)$  (see resp. Titchmarsh [10], ch. 7 and Heath-Brown [3]), inequality (2.10) gives

$$|\mathcal{S}_0^{(0)}| \ll \min(T^{2+\varepsilon} Y^{6-12\sigma}, T^{1+\varepsilon} Y^{2-4\sigma}),$$

hence the contribution of the zeros in  $\mathcal{S}_0^{(0)}$  to (2.4) is

$$(2.11) \quad \min \ll (T^{2(2k-1)+\varepsilon} Y^{(2k-1)(6-12\sigma)}, T^{2k-1+\varepsilon} Y^{(2k-1)(2-4\sigma)}).$$

## c) THE INDUCTIVE STEP

Let

$$\mathcal{N}_v^k = \{t: v \equiv n_k(t) \pmod{2v}\},$$

where  $v=2^l$  and  $l \ll \log T$ . From (2.9) we have

$$\begin{aligned} \sum_k &\ll \log T \sum_v \sum_{n_1, n_2} (n_1 n_2)^{-1/2} \left| \sum_{t \in \mathcal{N}_v^k} n_k(t) \left( \frac{n_2}{n_1} \right)^{it/2} \right|^2 \ll \\ &\ll \log T \sum_v \sum_{t_1, t_2 \in \mathcal{N}_v^k} n_k(t_1) n_k(t_2) \left| \sum_{n=M}^P n^{-1/2-i(t_1-t_2)} \right|^2 \ll \\ &\ll (\log T)^2 \max_v \left( v^2 \sum_{t_1, t_2 \in \mathcal{N}_v^k} \left| \sum_{n=M}^P n^{-1/2-i(t_1-t_2)} \right|^2 \right). \end{aligned}$$

Now we use Theorem 1 of Heath-Brown [4] (see also (9) of [2]), thus obtaining

$$\sum_k \ll T^\varepsilon \max_v v^2 (|\mathcal{N}_v^k| P + |\mathcal{N}_v^k|^2 + |\mathcal{N}_v^k|^{5/4} T^{1/2}).$$

But

$$\begin{aligned} v |\mathcal{N}_v^k| &\equiv \sum_t n_k(t) \ll N(\sigma, T)^k \\ v^2 |\mathcal{N}_v^k| &\equiv \sum_t n_k(t)^2 \ll R_{2k}, \end{aligned}$$

hence

$$(2.12) \quad \sum_k \ll T^\varepsilon (PR_{2k} + N(\sigma, T)^{2k} + R_{2k}^{3/4} N(\sigma, T)^{k/2} T^{1/2}).$$

From (2.7) and (2.12) we have

$$\begin{aligned} (2.13) \quad R_{2k} &\ll T^\varepsilon P^{1-2\sigma} (R_{2k} P + N(\sigma, T)^{2k} + R_{2k}^{3/4} N(\sigma, T)^{k/2} T^{1/2})^{1/2} \times \\ &\times (R_{2(k-1)} P + N(\sigma, T)^{2(k-1)} + R_{2(k-1)}^{3/4} N(\sigma, T)^{(k-1)/2} T^{1/2})^{1/2}. \end{aligned}$$

Now, as remarked in [4], if  $P \geq T^{2/3}$  the third term on the right of (2.12) may be neglected, so (2.13) gives

$$(2.14) \quad R_{2k} \ll T^\varepsilon P^{1-2\sigma} (R_{2k}^{1/2} P^{1/2} + N(\sigma, T)^k) (R_{2(k-1)}^{1/2} P^{1/2} + N(\sigma, T)^{k-1})$$

provided

$$(2.15) \quad P > T^{2/3}.$$

We write (2.14) as

$$(2.16) \quad R_{2k} \ll A_0 R_{2k}^{1/2} + A_1$$

where

$$\begin{aligned} A_0 &= T^\varepsilon P^{3/2-2\sigma} (R_{2(k-1)}^{1/2} P^{1/2} + N(\sigma, T)^{k-1}) \\ A_1 &= T^\varepsilon P^{1-2\sigma} N(\sigma, T)^k (R_{2(k-1)}^{1/2} P^{1/2} + N(\sigma, T)^{k-1}). \end{aligned}$$

From (2.16) we have either

$$R_{2k} \ll A_1$$

or

$$R_{2k} \ll A_0 R_{2k}^{1/2}, \quad \text{i.e.} \quad R_{2k} \ll A_0^2.$$

Hence from (2.16) we get

$$(2.17) \quad R_{2k} \ll A_0^2 + A_1$$

that is

$$(2.18) \quad R_{2k} \ll T^\varepsilon (R_{2(k-1)} P^{4(1-\sigma)} + R_{2(k-1)}^{1/2} P^{3/2-2\sigma} N(\sigma, T)^k + P^{3-4\sigma} N(\sigma, T)^{2(k-1)} + P^{1-2\sigma} N(\sigma, T)^{2k-1}).$$

Now we have either

$$(2.19) \quad R_{2(k-1)}^{1/2} P^{3/2-2\sigma} N(\sigma, T)^k \ll R_{2(k-1)} P^{4(1-\sigma)}$$

or

$$(2.20) \quad R_{2(k-1)} \ll P^{4\sigma-5} N(\sigma, T)^{2k}.$$

Suppose that (2.20) holds; then (2.18) gives

$$(2.21) \quad R_{2k} \ll T^\varepsilon (N(\sigma, T)^{2k} P^{-1} + P^{3-4\sigma} N(\sigma, T)^{2(k-1)} + P^{1-2\sigma} N(\sigma, T)^{2k-1}).$$

Suppose further that

$$(2.22) \quad N(\sigma, T) \gg P^{2(1-\sigma)}$$

in the sense that  $P^{2(1-\sigma)} \ll$  (estimate for  $N(\sigma, T)$ , cf. (2.29)). It is then easy to see that if (2.15), (2.20) and (2.22) hold we have

$$(2.23) \quad R_{2k} \ll T^\varepsilon N(\sigma, T)^{2k} P^{-1}.$$

Suppose now that (2.19) holds; then (2.18) gives

$$(2.24) \quad R_{2k} \ll T^\varepsilon (R_{2(k-1)} P^{4(1-\sigma)} + P^{1-2\sigma} N(\sigma, T)^{2k-1}),$$

provided (2.22) holds. We prove by induction on  $k$  that

$$(2.25) \quad R_{2k} \ll T^\varepsilon (P^{4(k-1)(1-\sigma)} R_2 + P^{1-2\sigma} N(\sigma, T)^{2k-1}).$$

In fact (2.24) is of the form

$$(2.26) \quad R_{2k} \ll B_0 R_{2(k-1)} + B_1(k)$$

where

$$B_0 = T^\varepsilon P^{4(1-\sigma)}, \quad B_1(k) = T^\varepsilon P^{1-2\sigma} N(\sigma, T)^{2k-1},$$

and from (2.26) it follows that

$$(2.27) \quad R_{2k} \ll B_0^{k-1} R_2 + B_1(k)$$

which is (2.25).

Indeed (2.27) is trivial for  $k=1$  and, assuming (2.27) to be true for  $k-1$ , we have from (2.26)

$$R_{2k} \ll B_0^{k-1} R_2 + B_1(k) + B_0 B_1(k-1);$$

but from (2.22) we see that  $B_0 B_1(k-1) \ll B_1(k)$ , so (2.27) holds. Since  $R_2 \ll N(\sigma, T)$  and, provided (2.22) holds,  $P^{1-2\sigma} N(\sigma, T)^{2k-1} \ll N(\sigma, T)^{2k} P^{-1}$ , we can summarize our results as

$$(2.28) \quad R_{2k} \ll T^\varepsilon (N(\sigma, T)^{2k} P^{-1} + P^{4(k-1)(1-\sigma)} N(\sigma, T)),$$

provided (2.15) and (2.22) hold.

d) THE FINAL ESTIMATE

Now we have to choose the parameters. Let

$$(2.29) \quad x = T^\varepsilon, \quad y = T^\theta, \quad N = T^{\theta'}, \quad N(\sigma, T) \ll T^{\beta(\sigma) + \varepsilon}$$

where  $\theta = \theta(\sigma)$ ,  $\theta' = \theta'(\sigma)$ ,  $0 < \theta' \leq \theta$ . Hence  $T^{\theta'h} \ll P \ll T^{\theta'h}$ ,  $h = h(\sigma)$ ,  $h \in \mathbb{N}$ ; let further  $\alpha(\sigma) = \theta'h$ . The optimal choice for  $\alpha(\sigma)$  would be

$$\alpha(\sigma) = \frac{(2k-1)\beta(\sigma)}{1+4(k-1)(1-\sigma)},$$

we choose

$$h = \left\lceil \frac{\beta(\sigma)}{\theta'} \cdot \frac{2k-1}{1+4(k-1)(1-\sigma)} \right\rceil,$$

hence

$$(2.30) \quad \alpha(\sigma) = \beta(\sigma) \frac{2k-1}{1+4(k-1)(1-\sigma)} - \xi\theta'$$

where  $0 \leq \xi \leq 1$ .

With the above choice of  $\alpha$  we get from (2.28):

$$(2.31) \quad R_{2k} \ll T^{\beta(\sigma) \left( 2k - \frac{2k-1}{1+4(k-1)(1-\sigma)} \right) + \xi\theta' + \varepsilon}.$$

Inserting (2.11) and (2.31) in (2.4) we obtain

$$(2.32) \quad N_{2k}(\sigma, T) \ll T^{\beta(\sigma) \left( 2k - \frac{2k-1}{1+4(k-1)(1-\sigma)} \right) + \theta(\sigma) + \varepsilon} + \min(T^{2(2k-1) + 6\theta(\sigma)(2k-1)(1-2\sigma) + \varepsilon}, T^{2k-1 + 2\theta(\sigma)(2k-1)(1-2\sigma) + \varepsilon}) = A + \min(B, G).$$

Equating  $A$  and  $B$ ,  $A$  and  $C$  we get respectively

$$(2.33) \quad \theta(\sigma) = \theta_B(\sigma) = \frac{2-\beta(\sigma)}{6(2\sigma-1)} + O\left(\frac{1}{k}\right)$$

$$(2.34) \quad \theta(\sigma) = \theta_G(\sigma) = \frac{1-\beta(\sigma)}{2(2\sigma-1)} + O\left(\frac{1}{k}\right).$$

We will use the following density estimates:

$$(2.35) \quad \beta(\sigma) = 2(1-\sigma), \quad \frac{11}{14} \leq \sigma \leq 1 \quad (\text{Jutila [7]})$$

$$(2.36) \quad \beta(\sigma) = \frac{1-\sigma}{7\sigma-5}, \quad \frac{43}{55} \leq \sigma \leq \frac{11}{14} \quad (\text{Jutila [7]})$$

$$(2.37) \quad \beta(\sigma) = \frac{9(1-\sigma)}{8\sigma-2}, \quad \frac{10}{13} \leq \sigma \leq \frac{43}{55} \quad (\text{Ivić [6]})$$

$$(2.38) \quad \beta(\sigma) = \frac{3(1-\sigma)}{7\sigma-4}, \quad \frac{3}{4} \leq \sigma \leq \frac{10}{13} \quad (\text{Ivić [6]})$$

$$(2.39) \quad \beta(\sigma) = \frac{3(1-\sigma)}{2-\sigma}, \quad \frac{1}{2} \leq \sigma \leq \frac{3}{4} \quad (\text{see Montgomery [8]}).$$

Since  $\beta(\sigma)$  is a decreasing function, in view of (2.38) we make the following choice of  $\theta(\sigma)$  in (2.32):

$$(2.40) \quad \theta(\sigma) = \theta_B(\sigma) \quad \text{if} \quad 10/13 \leq \sigma \leq 1$$

$$(2.41) \quad \theta(\sigma) = \theta_G(\sigma) \quad \text{if} \quad 1/2 \leq \sigma \leq 10/13.$$

Finally, we have to assure the compatibility of (2.30) with (2.15) and (2.22). It is easy to see that (2.30) and (2.22) are compatible for every  $\sigma \in (1/2, 1)$ .

Now let  $10/13 \leq \sigma \leq 1$ ,  $\beta(\sigma) = A(\sigma)(1-\sigma)$ ; we have to verify that

$$\frac{\beta(\sigma)}{2(1-\sigma)} > \frac{2}{3} + \frac{2-\beta(\sigma)}{6(2\sigma-1)} + O\left(\frac{1}{k}\right),$$

i.e. that

$$(2.42) \quad A(\sigma) > 2 + \frac{2(1-\sigma)}{5\sigma-2} + O\left(\frac{1}{k}\right).$$

Inserting (2.35)–(2.37) in (2.42) we see that (2.42) is never satisfied in the range  $10/13 \leq \sigma \leq 1$ .

If  $1/2 \leq \sigma \leq 10/13$  we have to verify that

$$(2.43) \quad A(\sigma) > \frac{8\sigma-1}{3\sigma} + O\left(\frac{1}{k}\right).$$

Inserting (2.38) and (2.39) we see that (2.43) is satisfied for

$$(2.44) \quad \sigma \in \left( \frac{1}{2} + O\left(\frac{1}{k}\right), \frac{6+\sqrt{22}}{14} + O\left(\frac{1}{k}\right) \right).$$

Summarizing our results we may write

$$(2.45) \quad N_{2k}(\sigma, T) \ll T^{\beta(\sigma)\left(2k - \frac{2k-1}{1+4(k-1)(1-\sigma)}\right) + \frac{1-\beta(\sigma)}{2(2\sigma-1)} + O\left(\frac{1}{k}\right) + \epsilon}$$

if  $\sigma \in \left(1/2 + O\left(\frac{1}{k}\right), \frac{6 + \sqrt{22}}{14} + O\left(\frac{1}{k}\right)\right)$ , where  $\beta(\sigma)$  is given by (2.38) and (2.39).

Since the trivial estimate is

$$N_{2k}(\sigma, T) \ll N(\sigma, T)^{2k-1} \log T,$$

(2.45) is non-trivial only when

$$(2.46) \quad \frac{\beta(\sigma)}{2(1-\sigma)} - \frac{1-\beta(\sigma)}{2(2\sigma-1)} > \beta(\sigma) + O\left(\frac{1}{k}\right).$$

Inserting (2.38) and (2.39) in (2.46) we see that (2.46) is satisfied only when

$$(2.47) \quad \sigma \in \left(\frac{2}{3} + O\left(\frac{1}{k}\right), \frac{6 + \sqrt{22}}{14} + O\left(\frac{1}{k}\right)\right).$$

The theorem is finally obtained from (2.38), (2.39), (2.45) and (2.47). ■

### 3. Appendix

In [9] we proved the asymptotic formula for  $\sum_{p \leq x} (\pi(p+h) - \pi(p))^{2k-1}$ , for each fixed positive integer  $k$ , under the assumption of the Riemann Hypothesis, where  $h \cong f_k(x)$  and  $f_k(x)$  is a suitable function of  $x$ , depending on  $k$ . We show here how such results may be obtained by a short direct argument. We have indeed the following

**THEOREM 2.** *Let  $k$  be a positive integer and let*

$$(3.1) \quad \psi(n+h) - \psi(n) \sim h \quad \text{for almost all } n,$$

with  $h \gg x^\theta$ ,  $0 < \theta < 1$ . Then

$$(3.2) \quad \sum_{p \leq x} (\psi(p+h) - \psi(p))^k \sim \frac{h^k x}{\log x}$$

for  $h \gg x^\theta$ .

**PROOF.** As in [9] we have

$$(3.3) \quad \int_0^x (\psi(t+h) - \psi(t))^{k+1} dt = \sum_{\substack{n_1, \dots, n_{k+1} \\ n \leq x, 0 \leq N-n \leq h}} \Lambda(n_1) \dots \Lambda(n_{k+1}) (h - N + n) +$$

$$+ \sum_{\substack{n_1, \dots, n_{k+1} \\ x < n \leq x+h, N-n \leq h}} \Lambda(n_1) \dots \Lambda(n_{k+1}) (x + h - N) = \sum_1 + \sum_2$$

where  $n = \min(n_1, \dots, n_{k+1})$ ,  $N = \max(n_1, \dots, n_{k+1})$ . But

$$(3.4) \quad \sum_1 = \int_0^h \psi_{k+1}(x, u) du$$

and

$$(3.5) \quad \sum_2 = O(h^{k+1} \log^k x)$$

where

$$\psi_{k+1}(x, u) = \sum_{\substack{n_1, \dots, n_{k+1} \\ n \leq x \\ 0 \leq N-n \leq n}} \Lambda(n_1) \dots \Lambda(n_{k+1}).$$

We estimate the integral in (3.3) in the following way: first we note that

$$(3.6) \quad \int_0^x (\psi(t+h) - \psi(t))^{k+1} dt = \sum_{n \leq x} (\psi(n+h) - \psi(n))^{k+1} + O(h^{k+1})$$

(we may suppose without loss of generality that  $h \in \mathbb{N}$ ) and then we use (3.1) in the sum (3.6) for almost all  $n$ , and Brun—Titchmarsh Theorem, i.e.  $\psi(n+h) - \psi(n) \ll \ll h \frac{\log n}{\log h}$ , for the remaining  $o(x)$  values of  $n$ . Thus we get

$$(3.7) \quad \int_0^x (\psi(t+h) - \psi(t))^{k+1} dt \sim x h^{k+1}.$$

From (3.3)—(3.7), using a well-known Tauberian argument, we obtain

$$(3.8) \quad \psi_{k+1}(x, h) \sim (k+1) x h^k.$$

Finally we have

$$(3.9) \quad \begin{aligned} & \sum_{m \leq x} \Lambda(m) (\psi(m+h) - \psi(m))^k = \\ & = \sum_{m \leq x} \Lambda(m) \left( \sum_{\substack{m_1, \dots, m_k \\ m \leq m_j < m+h}} \Lambda(m_1) \dots \Lambda(m_k) \right) \end{aligned}$$

and rearranging the sums in (3.9) we get, recalling the Remark a) of [9],

$$(3.10) \quad \sum_{m \leq x} \Lambda(m) (\psi(m+h) - \psi(m))^k \sim \frac{1}{k+1} \psi_{k+1}(x, h).$$

Theorem 2 follows now from (3.8) and (3.10) by partial summation.

**COROLLARY.** *Under the same assumptions of Theorem 2 we have  $\psi(p+h) - \psi(p) \sim h$  for almost all primes  $p$ , for  $h \gg x^\theta$ .*

**PROOF.** Using the estimate of Theorem 2 with  $k=1, 2$  we obtain

$$(3.11) \quad \sum_{x < p \leq 2x} (\psi(p+h) - \psi(p) - h)^2 = o\left(\frac{h^2 x}{\log x}\right).$$

Hence the cardinality of the set  $\mathcal{P}_0$  of primes  $p$  in  $(x, 2x]$  for which

$$\psi(p+h) - \psi(p) - h = o(h)$$



does not hold satisfies

$$|\mathcal{P}_0| \ll \frac{g_1(x) \frac{h^2 x}{\log x}}{g_2(x) h^2},$$

and the Corollary follows choosing in a suitable way the functions  $g_1(x)$  and  $g_2(x)$ , which tend to 0. ■

For instance, it follows from the Density Hypothesis that we may choose in Theorem 2  $\theta = \varepsilon$ , for every  $\varepsilon > 0$ , and from Huxley's density estimate in [5] we have the unconditional choice  $\theta = (1/6) + \varepsilon$ .

## REFERENCES

- [1] HEATH-BROWN, D. R., The difference between consecutive primes, II., *J. London Math. Soc.* **19** (1979), 207—220. *MR 80k*: 10041.
- [2] HEATH-BROWN, D. R., Zero density estimates for the Riemann Zeta-function and Dirichlet  $L$ -functions, *J. London Math. Soc.* **19** (1979), 221—232. *MR 80i*: 10055.
- [3] HEATH-BROWN, D. R., The twelfth power moment of the Riemann Zeta-function, *Quart. J. Math. Oxford Ser. (2)* **29** (1978), 443—462. *MR 80d*: 10059.
- [4] HEATH-BROWN, D. R., A large values estimate for Dirichlet polynomials, *J. London Math. Soc.* **20** (1979), 8—18. *MR 81a*: 10052.
- [5] HUXLEY, M. N., On the difference between consecutive primes, *Invent. Math.* **15** (1972), 164—170. *MR 45*# 1856.
- [6] Ivić, A., A Zero-density theorem for the Riemann Zeta-function (to appear).
- [7] JUTILA, M., Zero-density estimates for  $L$ -functions, *Acta Arith.* **32** (1977), 55—62. *MR 55*# 2800.
- [8] MONTGOMERY, H. L., *Topics in Multiplicative Number Theory*, Springer-Verlag, Berlin, 1971. *MR 49*# 2616.
- [9] PERELLI, A. and SALERNO, S., On an average of primes in short intervals, *Acta Arith.* **42** (1982), 91—96.
- [10] TITCHMARSH, E. C., *The Theory of the Riemann Zeta-function*, Oxford, 1951. *MR 13*—741.

(Received November 2, 1983)

ISTITUTO DI MATEMATICA  
UNIVERSITA' DI GENOVA  
VIA L. B. ALBERTI 4  
I-16132 GENOVA

ISTITUTO DI MATEMATICA  
FACOLTA' DI SCIENZE  
UNIVERSITA' DI SALERNO  
I-84100 SALERNO  
ITALY



# ON THE DISTRIBUTION OF $x_1^k + \dots + x_s^k$ IN THE ARITHMETIC PROGRESSIONS

SAVERIO SALERNO

1. One of the central problems in analytic number theory is the distribution of a sequence  $\mathcal{A} = \{a_n\}_{n \in \mathbb{N}}$  in arithmetic progressions. To make precise this problem, we define

$$(1.1) \quad A(x) = \sum_{n \leq x} a_n$$

$$(1.2) \quad A(x, d) = \sum_{\substack{n \leq x \\ n \equiv 0 \pmod{d}}} a_n.$$

Then, our purpose is to show that

$$(1.3) \quad A(x, d) = F(d)A(x) + R(x, d)$$

with  $F(d)$  multiplicative, where  $F(d)A(x)$  is the main term of  $A(x, d)$  and  $R(x, d)$  is the error, which has to be small "in average", that is

$$(1.4) \quad \sum_{d \leq x^\alpha} |R(x, d)| \ll \frac{A(x)}{(\log x)^N} \quad \forall N, \varepsilon > 0$$

for some  $\alpha$  with  $0 \leq \alpha \leq 1$ ; this number  $\alpha$  measures the level of distribution of  $\mathcal{A}$ . The behaviour of primes in arithmetic progressions is a typical example, and here Bombieri's theorem gives a satisfactory estimate for the error in average ( $\alpha = 1/2$  in (1.4)), whilst the conjecture of Halberstam and Richert would give the optimal level of distribution  $\alpha = 1$ .

Also, this problem is important not only in itself, but in connection with sieve methods, which are able to show the existence of  $r$ -almost primes  $P_r$  (that is, numbers with at most  $r$  prime factors, counting the multiplicity) in  $\mathcal{A}$ , (that is, for which  $a_n > C > 0$ ), if (1.3), (1.4) are known, with a value of  $r$  decreasing with respect to  $\alpha$  of (1.4). In this manner, one can approximate the problem of the representation of primes by  $\mathcal{A}$ , also if it is not possible at present to prove representation of primes by sieve methods, as it is shown by an example of Selberg, due to parity phenomenon (see later, formulas (1.11), (1.12),  $\delta_k$  can really assume every value between 0 and 2, an also can be oscillating). For an account to sieve methods, we refer to [5], [8], [11].

Bombieri [2] showed recently, that if (1.4) is known with  $\alpha=1$  and some other assumptions on  $F(d)$ ,  $R(x, d)$  are satisfied, then for every  $G \in C_0(T_r)$  and  $r \geq 2$

$$(1.5) \quad \sum_{\substack{n \leq x \\ n \in \mathcal{P}_r}} a_n G^*(n) \sim \left( \int_{T_r} G(u) \gamma_x(u) d\mu_r \right) \frac{HA(x)}{\log x}$$

where

$$(1.6) \quad \mathcal{P}_r = \{n = p_1 \dots p_r \text{ square-free}\},$$

$$(1.7) \quad T_r = \{(u_1, \dots, u_r) \in [0, 1]^r \mid u_1 + \dots + u_r = 1\}$$

$$(1.8) \quad d\mu_r = \frac{du_1 \dots du_{r-1}}{u_1 \dots u_r} \quad \text{on } T_r \quad \text{for } r \geq 2$$

$$(1.9) \quad G^*(n) = G^*(p_1 \dots p_r) = G\left(\frac{\log p_1}{\log n}, \dots, \frac{\log p_r}{\log n}\right)$$

$$(1.10) \quad H = - \sum_d \mu(d) \log d F(d) = \prod_p \frac{p}{p-1} (1 - F(p))$$

and, if we define  $\delta_x$  by means of

$$(1.11) \quad \sum_{p \leq x} a_p \sim \delta_x \frac{HA(x)}{\log x}$$

we have

$$(1.12) \quad \gamma_x(u) = \begin{cases} \delta_x & \text{if } r \text{ is odd} \\ 2 - \delta_x & \text{if } r \text{ is even.} \end{cases}$$

Here,  $\delta_x$  and  $\gamma_x(u)$  are naturally defined only mod  $o(1)$ .

Of course, condition (1.4) with  $\alpha=1$ , is in general very difficult to prove in the applications of Bombieri's asymptotic sieve, in which  $A(x) \gg x^{1-\varepsilon}$ , as pointed out by Bombieri itself. Moreover, if  $A(x) \sim x^{\theta}$  with  $\theta < 1$ , it is probably also false; consider for instance the case  $\mathcal{A} = \{n^2 + 1\}_{n \in \mathbb{N}}$ ; then  $A(x) \sim x^{1/2}$ ,  $F(d) \sim 1/d$  and, for  $d > x^{(1/2)+\varepsilon}$ , we have  $F(d)A(x) = O(x^{-\varepsilon}) = o(1)$ , whilst  $A(x, d)$  is a non-negative integer number, so we have

$$(1.13) \quad |R(x, d)| = |A(x, d) - F(d)A(x)| > \frac{A(x, d)}{2}.$$

Hence (1.4) with  $\alpha > 1/2$  is not realistic. In this case, a non-trivial sharper treatment of the error term of the sieve is needed, and this is accomplished by the new bilinear form in Rosser's sieve given by Iwaniec [6], which among other things, enables him to show that  $n^2 + 1$  represents  $P_2$  [7]. Also for Selberg's sieve it is possible to put the error term in such bilinear form (see [9]—[10]).

Nevertheless, in the case  $A(x) \gg x^{1-\varepsilon}$ , condition (1.4) can be conjectured and the purpose of the present paper is to prove it for the sequence

$$(1.14) \quad \mathcal{A} = \{\mu^2(n)r(n)\}_{n \in \mathbb{N}, (n, K)=1}$$

where

$$(1.15) \quad r(n) = \# \text{ sol } \{x_1^k + \dots + x_s^k = n, x_i \geq 0\}$$

if  $s$  is so large with respect to  $k$  that the asymptotic formula for Waring problem

$$(1.16) \quad r(n) = \frac{\Gamma^s\left(1 + \frac{1}{k}\right)}{\Gamma\left(\frac{s}{k}\right)} \mathfrak{S}(n) n^{(s/k)-1} + O(n^{(s/k)-1-\delta}) \quad \text{for some } \delta > 0$$

is valid; here

$$(1.17) \quad \mathfrak{S}(n) = \sum_{q=1}^{\infty} A(q, n)$$

$$(1.18) \quad A(q, n) = \frac{1}{q^s} \sum_{a \in \mathbb{Z}_q^*} S^s(a, q) e^{-2\pi i a n / q}$$

$$(1.19) \quad S(a, q) = \sum_{t=1}^q e^{2\pi i a t^k / q}$$

and we have

$$(1.20) \quad 1 \ll \mathfrak{S}(n) \ll 1.$$

The asymptotic formula for  $r(n)$  has been proved by Vinogradov [13] using his method of trigonometrical sums for

$$(1.21) \quad s > (4 + \varepsilon) k^2 \log k$$

if  $k \gg k_0$ , improving on the earlier values of  $s(k) \sim k 2^{k-1}$  due to Hardy and Littlewood. For a survey of Waring's problem, we also refer to Vaughan [12] and Davenport [3].

In this paper, we shall always assume, also without explicit mention, that  $s$  is large enough to ensure (1.16).

I take the pleasure to thank Professor Bombieri for the helpful discussions on the subject, and Professor Pintz for constant encouragement and suggestions.

2. In this section, we state our main results. As it will be seen, they require some assumption on  $k$ , in order to avoid further technical difficulties. Moreover, for the same reason we shall restrict our sequences to the integers coprimes with  $k$ .

**THEOREM 1.** *Let  $\mathfrak{S}(n)$  be defined by (1.17),  $k$  be square-free. We define*

$$S_d = \sum_{\substack{n < x, (n, k) = 1 \\ n \equiv 0 \pmod{d}}} \mu^2(n) \mathfrak{S}(n).$$

*Then we have*

$$(2.1) \quad S_d = F(d) S_1 + R_d$$

*where*

$$(2.2) \quad S_1 = Lx + O(x(\log x)^{-s/3k})$$

*with  $L$  defined by (5.21), and  $F(d)$  is a multiplicative function defined by*

$$(2.3) \quad F(d) = \frac{d}{h(d)} \prod_{p|d} \left(1 - \frac{1}{p^2}\right) \left\{ \sum_{t|d} \frac{g(t)}{h(t)} \right\} \frac{f(d) M(d, 0)}{d^{s-1}} \prod_{p|d} \left(1 + \frac{1}{p}\right)^{-1} \frac{1}{d}$$

with  $f(d)$ ,  $g(d)$ ,  $h(d)$ ,  $M(d, 0)$  defined by

$$(2.4) \quad M(d, a) = \# \text{ sol } \{x_1^k + \dots + x_s^k \equiv a \pmod{d}, x_i \in \mathbb{Z}_d\}$$

$$(2.5) \quad f(d) = \prod_{\substack{p|d \\ p \nmid k}} \left(1 - \frac{1}{M(p, 0)}\right)$$

$$(2.6) \quad g(d) = \varphi(d) \frac{f(d)}{d} \frac{M(d, 0)}{d^{s-1}}$$

$$(2.7) \quad h(d) = \sum_{e \in \mathbb{Z}_d^*} \frac{M(d, e)}{d^{s-1}}.$$

Moreover, for the error we have

$$(2.8) \quad R_d \ll \frac{x}{d(\log x)^{s/3k}} + x^\delta \sqrt{\frac{x}{d}} \quad \forall \delta > 0$$

$$(2.9) \quad \sum_{d < x^{1-\delta}} |R_d| \ll \frac{S_1}{(\log x)^{(s/3k)-1}}.$$

Finally,  $S_1$  has a distribution function on the primes  $\delta_x = \delta$  (constant) defined by

$$(2.10) \quad \delta = \frac{k^2}{\varphi(k^2)} \left\{ \prod_{p \nmid k} \left(1 + \frac{g(p)}{h(p)}\right) (1 - F(p)) \right\}^{-1}, \quad (0 < \delta < 2).$$

**THEOREM 2.** Let  $k$  be square-free,  $r(n)$  be defined by (1.15). Suppose  $s > s_0(k)$  large enough to ensure the asymptotic formula (1.16).

We define

$$N_d = \sum_{\substack{n < x, (n, k) = 1 \\ n \equiv 0 \pmod{d}}} \mu^2(n) r(n).$$

Then we have

$$(2.11) \quad N_d = F(d) N_1 + \tilde{R}_d$$

where

$$(2.12) \quad N_1 = \frac{k}{s} L x^{s/k} + O(x^{s/k} (\log x)^{-s/3k})$$

with  $L$  given by (5.21) and  $F(d)$  the multiplicative function defined by (2.3). Moreover

$$(2.13) \quad \tilde{R}_d \ll \frac{x^{s/k}}{d(\log x)^{s/3k}} + x^{(s/k)-1+\delta} \sqrt{\frac{x}{d}} \quad \forall \delta > 0$$

$$(2.14) \quad \sum_{d < x^{1-\delta}} |\tilde{R}_d| \ll N_1 (\log x)^{-(s/3k)+1}.$$

Finally,  $N_1$  has the same distribution function on the primes  $\delta_x = \delta$  given by (2.10).

3. The first step in the proof of our results is to obtain a more handable expression for  $\mathfrak{S}(n)$ . This is accomplished by means of the following lemmas:

LEMMA 3.1. Suppose  $p \nmid k$ ,  $2 \leq \alpha \leq k$ . Then

$$(3.1) \quad S(a, p^\alpha) = p^{\alpha-1}.$$

PROOF. This is Lemma 4 of ch. 2 of [13]. ■

LEMMA 3.2. Suppose  $k$  square-free,  $p \nmid k$ ,  $3 \leq \alpha \leq k$ . Then

$$S(a, p^\alpha) = p^{\alpha-1}.$$

PROOF. We have by definition

$$(3.2) \quad S(a, p^\alpha) = \sum_{z=0}^{p^\alpha-1} e^{2\pi i a z^k / p^\alpha}.$$

We set

$$(3.3) \quad z = t p^{\alpha-2} + b \quad \text{with} \quad 0 \leq b \leq p^{\alpha-2} - 1, \quad 0 \leq t \leq p^2 - 1$$

and we remark that  $\alpha - 2 \geq 1$  since  $\alpha \geq 3$ . We have

$$(3.4) \quad S(a, p^\alpha) = \sum_{b=0}^{p^{\alpha-2}-1} e^{2\pi i a b^k / p^\alpha} \sum_{t=0}^{p^2-1} e^{2\pi i k t b^{k-1} / p^3}$$

since  $z^k \equiv b^k + k t p^{\alpha-2} b^{k-1} \pmod{p^\alpha}$  because  $p \nmid k$ .

Using the fact that  $k$  is square-free, we set

$$t = p t_1 + r \quad \text{with} \quad 0 \leq t_1 \leq p-1, \quad 0 \leq r \leq p-1$$

$$k = p k_1, \quad p \nmid k_1$$

and we have, using  $(a, p) = 1$

$$(3.5) \quad \begin{aligned} \sum_{t=0}^{p^2-1} e^{2\pi i a k b^{k-1} t / p^3} &= \sum_{t_1=0}^{p-1} \sum_{r=0}^{p-1} e^{2\pi i a k_1 b^{k-1} r / p} = \\ &= p \sum_{r=0}^{p-1} e^{2\pi i a k_1 b^{k-1} r / p} = \begin{cases} 0 & \text{if } p \nmid b \\ p^2 & \text{if } p \mid b. \end{cases} \end{aligned}$$

By (3.4), (3.5) we get

$$S(a, p^\alpha) = p^2 \sum_{\substack{b=0 \\ p \mid b}}^{p^{\alpha-2}-1} e^{2\pi i a b^k / p^\alpha} = p^2 p^{\alpha-3} = p^{\alpha-1}. \quad \blacksquare$$

LEMMA 3.3. Let  $\alpha \geq k+1$ . Then

$$(3.6) \quad S(a, p^\alpha) = p^{k-1} S(a, p^{\alpha-k}).$$

PROOF. This is Lemma 5 of ch. 2 of [13]. ■



LEMMA 3.4. *Let  $p \nmid k, n$  be square-free. Then*

$$(3.7) \quad A(p^\alpha, n) = \begin{cases} 0 & \text{if } \alpha \equiv 3 \text{ or if } \alpha = 2 \text{ and } p \nmid n \\ -p^{1-s} & \text{if } \alpha = 2 \text{ and } p \mid n \\ A(p, n) & \text{if } \alpha = 1. \end{cases}$$

PROOF. Let us write

$$(3.8) \quad \alpha = \beta k + \gamma \quad \text{with} \quad 0 \leq \gamma \leq k-1$$

and suppose first  $\gamma \neq 1$ . Then, using Lemmas 1.3 we have

$$(3.9) \quad A(p^\alpha, n) = \frac{1}{p^{\alpha s}} \sum_{a \in \mathbb{Z}_{p^\alpha}^*} \{S(a, p^\alpha)\}^s e^{-2\pi i a n / p^\alpha} = \frac{1}{p^{s(\beta+1)}} \sum_{a \in \mathbb{Z}_{p^\alpha}^*} e^{-2\pi i a n / p^\alpha}.$$

Now, if  $p \nmid n$ , we have  $a \in \mathbb{Z}_{p^\alpha}^* \Leftrightarrow a n \in \mathbb{Z}_{p^\alpha}^*$ ; hence, if  $p \nmid n$

$$(3.10) \quad A(p^\alpha, n) = \frac{1}{p^{\alpha s}} \sum_{a \in \mathbb{Z}_{p^\alpha}^*} e^{-2\pi i a n / p^\alpha} = \frac{1}{p^{\alpha s}} \mu(p^\alpha) = 0$$

since  $\alpha > 1$ . If  $p \mid n$ , then  $p^2 \nmid n$  since  $n$  is square-free and we set  $n = p n_1$ ,  $p \nmid n_1$ ,

$$a = \mu p^{\alpha-1} + v \quad \text{with} \quad v \in \mathbb{Z}_{p^{\alpha-1}}^* \quad \text{and} \quad 0 \leq \mu \leq p-1$$

obtaining, since  $p \nmid n_1$

$$(3.11) \quad \begin{aligned} A(p^\alpha, n) &= \frac{1}{p^{s(\beta+1)}} \sum_{a \in \mathbb{Z}_{p^\alpha}^*} e^{-2\pi i a n_1 / p^{\alpha-1}} = \frac{1}{p^{s(\beta+1)}} \sum_{\substack{0 \leq \mu \leq p-1 \\ v \in \mathbb{Z}_{p^{\alpha-1}}^*}} e^{-2\pi i \mu n_1 / p^{\alpha-1}} = \\ &= \frac{p}{p^{s(\beta+1)}} \sum_{v \in \mathbb{Z}_{p^{\alpha-1}}^*} e^{-2\pi i v / p^{\alpha-1}} = \frac{p}{p^{s(\beta+1)}} \mu(p^{\alpha-1}). \end{aligned}$$

Hence, if  $\alpha \not\equiv 1 \pmod{k}$  we get if  $p \nmid n$

$$(3.12) \quad A(p^\alpha, n) = \begin{cases} 0 & \text{if } \alpha \equiv 3 \\ -\frac{p}{p^{s(\beta+1)}} & \text{if } \alpha = 2. \end{cases}$$

Now, we come to the case  $\alpha = \beta k + 1$ . We have using Lemma 3

$$(3.13) \quad A(p^\alpha, n) = \frac{1}{p^{s(\beta+1)}} \sum_{a \in \mathbb{Z}_{p^\alpha}^*} \{S(a, p)\}^s e^{-2\pi i a n / p^\alpha}.$$

Setting

$$a = t p + r \quad \text{with} \quad 1 \leq r \leq p-1, \quad 0 \leq t \leq p^{\alpha-1}-1$$

we have, since  $S(a, p) = S(r, p)$

$$(3.14) \quad A(p^\alpha, n) = \frac{1}{p^{s(\beta+1)}} \sum_{r=1}^{p-1} \{S(r, p)\}^s e^{-2\pi i r n / p^\alpha} \sum_{t=0}^{p^{\alpha-1}-1} e^{-2\pi i t n / p^{\alpha-1}} =$$

$$= \begin{cases} 0 & \text{if } p^{\alpha-1} \nmid n \\ \frac{p^{\alpha-1}}{p^{s(\beta+1)}} A(p, n) & \text{if } p^{\alpha-1} \mid n. \end{cases}$$

Recalling that  $n$  is square-free, (3.14) becomes

$$(3.15) \quad A(p^\alpha, n) = \begin{cases} 0 & \text{if } \alpha > 1 \\ A(p, n) & \text{if } \alpha = 1 \end{cases} \quad \text{for } \alpha \equiv 1 \pmod{k}.$$

Now, our Lemma follows from (3.10), (3.12) and (3.15). ■

LEMMA 3.5. Let  $p \mid k$ ,  $k$  square-free; then

$$(3.16) \quad A(p^\alpha, n) = \begin{cases} 0 & \text{if } \alpha \equiv 3 \\ A(p^2, n) & \text{if } \alpha = 2 \\ A(p, n) & \text{if } \alpha = 1. \end{cases}$$

PROOF. Since  $k$  is square-free, we have  $k = pk_1$  with  $p \nmid k_1$ . We also write  $\alpha = \beta k + \gamma$  with  $0 \leq \gamma \leq k-1$ . If  $\gamma = 1$  or  $\gamma \equiv 3$ , we can now use Lemma 2 (in place of Lemma 1) and Lemma 3, and we obtain, as in Lemma 4, formulas (3.10), (3.12), (3.15). So, it remains only to examine the case  $\gamma = 2$ , for which we also use Lemmas 2 and 3. We get

$$(3.17) \quad \begin{aligned} A(p^\alpha, n) &= \frac{1}{p^{\alpha s}} \sum_{a \in \mathbb{Z}_{p^\alpha}^*} \{S(a, p^\alpha)\}^s e^{-2\pi i a n / p^\alpha} = \\ &= \frac{1}{p^{s(\beta+2)}} \sum_{a \in \mathbb{Z}_{p^2}^*} \{S(a, p^2)\}^s e^{-2\pi i a n / p^2}. \end{aligned}$$

Setting

$$a = tp^2 + r, \quad \text{with } 0 \leq t \leq p^{\alpha-2} - 1, \quad r \in \mathbb{Z}_{p^2}^*$$

we have, since  $S(a, p^2) = S(r, p^2)$

$$(3.18) \quad A(p^\alpha, n) = \frac{1}{p^{s(\beta+2)}} \sum_{r \in \mathbb{Z}_{p^2}^*} \{S(r, p^2)\}^s e^{-2\pi i r n / p^2} \sum_{t=0}^{p^{\alpha-2}-1} e^{-2\pi i t n / p^{\alpha-2}} =$$

$$= \begin{cases} 0 & \text{if } p^{\alpha-2} \nmid n \\ p^{\alpha-2-s(\beta+2)} \sum_{r \in \mathbb{Z}_{p^2}^*} \{S(r, p^2)\}^s e^{-2\pi i r n / p^2} & \text{if } p^{\alpha-2} \mid n. \end{cases}$$

Recalling that  $n$  is square-free, (3.18) becomes

$$(3.19) \quad A(p^2, n) = \begin{cases} 0 & \text{if } \alpha > 2 \\ A(p^2, n) & \text{if } \alpha = 2 \end{cases} \quad \text{for } \alpha \equiv 2 \pmod{k}$$

which, with the previous remark for  $\alpha \not\equiv 2 \pmod{k}$ , proves our Lemma. ■

We summarize the results of Lemma 4 and 5 in the following form.

LEMMA 3.6. *Let  $k, n$  be square-free. Then we have*

$$(3.20) \quad A(p^2, n) = \begin{cases} 0 & \text{if } \alpha \equiv 3 \text{ or if } \alpha = 2 \text{ and } p \nmid nk \\ -p^{1-s} & \text{if } \alpha = 2 \text{ and } p \mid n, p \nmid k \\ A(p^2, n) & \text{if } \alpha = 2 \text{ and } p \mid k \\ A(p, n) & \text{if } \alpha = 1. \end{cases}$$

Now, we introduce

$$(3.21) \quad M(q, n) = \# \text{ sol } \{x_1^k + \dots + x_s^k \equiv n \pmod{p}, x_i \in \mathbb{Z}_q\}$$

$$(3.22) \quad \psi(p, n) = \sum_{t=0}^{\infty} A(p^t, n).$$

From Vinogradov's book [13] we quote the following results:

$$(3.23) \quad r(n) = \frac{\Gamma\left(1 + \frac{1}{k}\right)^s}{\Gamma\left(\frac{s}{k}\right)} \mathfrak{S}(n) n^{(s/k)-1} + O(n^{(s/k)-1-(1/k^2)})$$

where

$$(3.24) \quad \mathfrak{S}(n) = \sum_{q=1}^{\infty} A(q, n);$$

this is the Theorem of ch. 7; moreover Lemma 11 of ch. 2:

$$(3.25) \quad \mathfrak{S}(n) = \prod_p \psi(p, n),$$

Lemma 10 of ch. 2:

$$(3.26) \quad \sum_{q \mid m} A(q, m) = \frac{M(m, n)}{m^{s-1}},$$

Lemmas 6 and 12 of ch. 2:

$$(3.27) \quad 1 \ll \mathfrak{S}(n) \ll 1.$$

Finally, by (3.22), (3.25), (3.26) and Lemma 6 we easily obtain

LEMMA 3.7. *Let  $k, n$  be square-free. Then we have*

$$(3.28) \quad \mathfrak{S}(n) = \prod_{p \mid k} \frac{M(p^2, n)}{p^{2(s-1)}} \prod_{p \nmid k} \frac{M(p, n)}{p^{s-1}} \prod_{\substack{p \mid n \\ p \nmid k}} \left(1 - \frac{1}{M(p, 0)}\right). \quad \blacksquare$$

4. In this section we give an asymptotic evaluation for

$$(4.1) \quad S_d = \sum_{\substack{n \leq x, (n, k)=1 \\ n \equiv 0 \pmod{d}}} \mu^2(n) \mathfrak{S}(n).$$

Due to condition  $(n, k)=1$ , we can assume  $(d, k)=1$ , because otherwise  $S_d$  is obviously zero, and also  $d$  is square-free, due to  $\mu^2(n)$ .

We need the following Lemma, for which we give an elementary proof:

LEMMA 4.1. *Let  $d$  be square-free,  $B$  an integer with  $(d, B)=1$ ,  $(a, B)=1$ . Then we have*

$$(4.2) \quad \sum_{\substack{n \leq x, (n, d)=1 \\ n \equiv a \pmod{B}}} \mu^2(n) = \frac{6}{\pi^2} \frac{1}{\varphi(B)} \prod_{p|dB} \left(1 + \frac{1}{p}\right)^{-1} x + O(B\tau(d)\sqrt{x})$$

where  $\tau(d)$  denotes the divisor-function.

PROOF. By a well-known formula, we have

$$(4.3) \quad \sum_{\substack{n \leq x, (n, d)=1 \\ n \equiv a \pmod{B}}} \mu^2(n) = \sum_{\substack{n \leq x \\ (n, d)=1}} \mu^2(n) \frac{1}{\varphi(B)} \sum_{\chi \pmod{B}} \bar{\chi}(a) \chi(n).$$

Now

$$(4.4) \quad \begin{aligned} \sum_{\substack{n \leq x \\ (n, d)=1}} \mu^2(n) \chi(n) &= \sum_{n \leq x} \chi(n) \sum_{t^2|n} \mu(t) \sum_{r|(n, d)} \mu(r) = \\ &= \sum_{r|d} \mu(r) \sum_{t \leq \sqrt{x}} \mu(t) \sum_{\substack{n \leq x \\ [r, t^2]|n}} \chi(n). \end{aligned}$$

Since

$$(4.5) \quad \sum_{\substack{n \leq x \\ [r, t^2]|n}} \chi(n) = \sum_{m \leq x/[r, t^2]} \chi([r, t^2]) \chi(m) = \begin{cases} \chi_0([r, t^2]) \sum_{\substack{m \leq x/[r, t^2] \\ (m, B)=1}} 1 & \text{if } \chi = \chi_0 \\ O(B) & \text{if } \chi \neq \chi_0 \end{cases}$$

where  $\chi_0$  denotes the principal character (mod  $B$ ), we obtain

$$(4.6) \quad \sum_{\substack{n \leq x \\ [r, t^2]|n}} \chi(n) = \chi_0([r, t^2]) \frac{\varphi(B)}{B} \frac{x}{[r, t^2]} E_0(\chi) + O(B)$$

where

$$E_0(\chi) = \begin{cases} 1 & \text{if } \chi = \chi_0 \\ 0 & \text{if } \chi \neq \chi_0. \end{cases}$$

By (4.6), formula (4.4) becomes

$$(4.7) \quad \sum_{\substack{n \leq x \\ (n, d)=1}} \mu^2(n) \chi(n) = \frac{\varphi(B)}{B} \left\{ \sum_{\substack{r|d \\ (r, B)=1}} \sum_{\substack{t \leq \sqrt{x} \\ (t, B)=1}} \frac{\mu(r)\mu(t)}{[r, t^2]} \right\} E_0(\chi) x + O(B\tau(d)\sqrt{x}).$$

Using  $(d, B)=1$ ,  $d$  square-free, we have

$$(4.8) \quad S = \sum_{\substack{r|d \\ (r, B)=1}} \sum_{\substack{t \leq \sqrt{x} \\ (t, B)=1}} \frac{\mu(r)\mu(t)}{[r, t^2]} = \sum_{\substack{t \leq \sqrt{x} \\ (t, B)=1}} \frac{\mu(t)}{t^2} \sum_{r|d} \frac{\mu(r)}{r} (r, t).$$

Since

$$(4.9) \quad g_t(d) = \sum_{r|d} \frac{\mu(r)}{r} (r, t) = \prod_{p|d} \left(1 - \frac{(p, t)}{p}\right) = \begin{cases} \prod_{p|d} \left(1 - \frac{1}{p}\right) & \text{if } (t, d) = 1 \\ 0 & \text{if } (t, d) > 1 \end{cases}$$

we get

$$(4.10) \quad \begin{aligned} S &= g(d) \sum_{\substack{t \leq \sqrt{x} \\ (t, B)=1}} \frac{\mu(t)}{t^2} = g(d) \sum_{\substack{t=1 \\ (t, dB)=1}}^{\infty} \frac{\mu(t)}{t^2} + O(x^{-1/2}) = \\ &= \frac{1}{\zeta(2)} \prod_{p|dB} \left(1 - \frac{1}{p^2}\right)^{-1} g(d) + O(x^{-1/2}). \end{aligned}$$

Now, our Lemma follows from (4.7), (4.8) and (4.10). ■

We introduce a number

$$(4.11) \quad A = k^2 \prod_{\substack{p < T \\ p \nmid k}} p$$

where  $T > k$  is a parameter to be chosen in the sequel.

We define

$$(4.12) \quad f(n) = \prod_{\substack{p < T \\ p|n, p \nmid k}} \left(1 - \frac{1}{M(p, 0)}\right).$$

Now, we are able to prove the following

LEMMA 4.2. For  $(d, k) = 1$ ,  $d, k$  square-free, we have for every  $\delta > 0$

$$(4.13) \quad \begin{aligned} S_d &= \frac{6}{\pi^2} \left\{ \sum_{\substack{v \in \mathbb{Z}_A \\ (v, k)=1 \\ v \equiv 0 \pmod{\alpha}}} f(v) \frac{M(A, v)}{A^{s-1}} \frac{\varphi((A_1, v))}{(A_1, v)} \right\} \frac{1}{\varphi(A_1)} \prod_{p|dA_1} \left(1 + \frac{1}{p}\right)^{-1} \frac{x}{d} \times \\ &\quad \times (1 + O(T^{-s/3})) + O(A^2 x^\delta \sqrt{x/d}) \end{aligned}$$

where  $A_1 = A/(A, d)$ ,  $\alpha = (A, d)$ .

PROOF. First of all, by a result of A. Weil [14] we have

$$(4.14) \quad M(p, n) = p^{s-1} + O(p^{s-1/2})$$

for which we easily obtain

$$(4.15) \quad \prod_{p > T} \frac{M(p, n)}{p^{s-1}} = \prod_{p > T} \{1 + O(p^{-(s-1)/2})\} = 1 + O(T^{-s/3})$$

$$(4.16) \quad \prod_{p > T} \left(1 - \frac{1}{M(p, 0)}\right) = \prod_{p > T} \{1 + O(p^{1-s})\} = 1 + O(T^{-s/2}).$$

Using the multiplicativity of  $M(q, n)$  in  $q$  and Lemma 3.7, we obtain by (4.7),

$$(4.17) \quad \mathfrak{S}(n) = f(n) \frac{M(A, n)}{A^{s-1}} + O\left(\frac{f(n)M(A, n)}{T^{s/3} A^{s-1}}\right).$$

Hence

$$(4.18) \quad S_d = \sum_{\substack{n \leq x \\ n \equiv 0 \pmod{d}}} \mu^2(n) f(n) \frac{M(A, n)}{A^{s-1}} \{1 + O(T^{-s/3})\}.$$

We have

$$(4.19) \quad \sum = \sum_{\substack{n \leq x \\ n \equiv 0 \pmod{d}}} \mu^2(n) f(n) \frac{M(A, n)}{A^{s-1}} = \sum_{\substack{v \in \mathbb{Z}_A \\ (v, k)=1}} f(v) \frac{M(A, v)}{A^{s-1}} \sum_{\substack{n_1 \leq x/d \\ (n_1, kd)=1 \\ n_1 d \equiv v(A)}} \mu^2(n_1)$$

since  $(n, k)=1, k|A$  and  $n \equiv v(A)$  imply  $(v, k)=1$ .

Now we set

$$(4.20) \quad \alpha = (d, A), \quad d = d_1 \alpha, \quad A = A_1 \alpha \quad \text{and} \quad (d, A_1) = 1$$

since  $(d, k)=1$  and  $d, A/k^2$  are square-free.

$$(4.21) \quad t = (v, A_1) \quad v = tv_1, \quad A_1 = tB \quad \text{and} \quad (v, B) = 1$$

since  $(v, k)=1$  and  $B/k^2$  is square-free; moreover,  $v \equiv O(\alpha)$  because  $nd \equiv O(\alpha)$ .

Then, we have

$$(4.22) \quad \sum_{\substack{n < x/d \\ (n, kd)=1 \\ nd \equiv v(A)}} \mu^2(n) = \sum_{\substack{n < x/d \\ (n, kd)=1 \\ n \equiv d^{-1}v(A_1)}} \mu^2(n) = \sum_{\substack{n < x/dt \\ (n, kdt)=1 \\ n \equiv d^{-1}v(B)}} \mu^2(n)$$

and in the last sum the condition  $(n, k)=1$  can be eliminated, since

$$n \equiv d^{-1}v(B), \quad k^2 | B \quad \text{and} \quad (d^{-1}v, B) = 1.$$

Using Lemma 4.1, we obtain since  $(d^{-1}v, B)=1, (dt, B)=1$ ,

$$(4.23) \quad \sum_{\substack{n < x/dt \\ (n, dt)=1 \\ n \equiv d^{-1}v(B)}} \mu^2(n) = \frac{6}{\pi^2} \frac{1}{\varphi(B)} \prod_{p|dtB} \left(1 + \frac{1}{p}\right)^{-1} \frac{x}{dt} + O\left(B\tau(dt) \left(\frac{x}{dt}\right)^{1/2}\right)$$

and clearly

$$(4.24) \quad \prod_{p|dtB} \left(1 + \frac{1}{p}\right)^{-1} = \prod_{p|dA_1} \left(1 + \frac{1}{p}\right)^{-1}.$$

Formula (4.23) produces in (4.19) an error which is

$$(4.25) \quad \ll x^\delta \sqrt{\frac{x}{d}} \sum_{\substack{v \in \mathbb{Z}_A \\ (v, k)=1 \\ v \equiv 0(\alpha)}} \frac{M(A, v)}{A^{s-1}} \frac{A_1}{(v, A_1)^{3/2}} \ll A^2 x^\delta \sqrt{\frac{x}{d}} \quad \forall \delta > 0.$$

The main term of (4.19) is

$$(4.26) \quad \Sigma \sim \frac{6}{\pi^2} \frac{1}{\varphi(A_1)} \prod_{p|dA_1} \left(1 + \frac{1}{p}\right)^{-1} \left\{ \sum_{\substack{v \in \mathbb{Z}_{A_1} \\ (v, k)=1 \\ v \equiv 0(a)}} f(v) \frac{M(A, v)}{A^{s-1}} \frac{\varphi((A_1, v))}{(A_1, v)} \right\} \frac{x}{d}.$$

Now, Lemma 4.2 follows from (4.18), (4.19), (4.25) and (4.26). ■

LEMMA 4.3. For  $(d, k)=1$ ,  $d, k$  square-free we have, for every  $\delta > 0$

$$(4.27) \quad S_d = F_T(d) S_1 + O\left(\frac{x}{d} T^{-s/3} + A^2 x^\delta \sqrt{\frac{x}{d}}\right)$$

where, for  $\alpha(d, A)$ ,  $F_T(d)$  is a multiplicative function defined as follows

$$(4.28) \quad F_T(d) = \frac{\alpha}{h(\alpha)} \prod_{p|\alpha} \left(1 - \frac{1}{p^2}\right) \left\{ \sum_{t|\alpha} \frac{g(t)}{h(t)} \right\}^{-1} \frac{f(\alpha) M(\alpha, 0)}{\alpha^{s-1}} \prod_{p|d} \left(1 + \frac{1}{p}\right)^{-1} \frac{1}{d}$$

with  $g(t)$ ,  $h(t)$  defined by (4.39), (4.40) of the sequel.

PROOF. We apply Lemma 4.2 for  $d=1$ ; we obtain

$$(4.29) \quad S_1 = \frac{6}{\pi^2} \frac{1}{\varphi(A)} \prod_{p|A} \left(1 + \frac{1}{p}\right)^{-1} \left\{ \sum_{\substack{v \in \mathbb{Z}_A \\ (v, k)=1}} f(v) \frac{M(A, v)}{A^{s-1}} \frac{\varphi((A, v))}{(A, v)} \right\} \times \\ \times (1 + O(T^{-s/3})) x + O(A^2 x^{1/2+\delta})$$

for every  $\delta > 0$ . Clearly

$$(4.30) \quad S_1 \ll \frac{1}{A} \sum_{v \in \mathbb{Z}_A} \frac{M(A, v)}{A^{s-1}} x \ll x.$$

Similarly we get, since  $(\alpha, A_1)=1$

$$(4.31) \quad S_d \ll \frac{x}{d} \frac{1}{A_1} \sum_{\substack{v \in \mathbb{Z}_{A_1} \\ v \equiv 0(\alpha)}} \frac{M(A, v)}{A^{s-1}} \ll \frac{M(\alpha, 0)}{\alpha^{s-1}} \frac{1}{A_1^s} \sum_{v \in \mathbb{Z}_{A_1}} M(A, \alpha v) \frac{x}{d} \ll \\ \ll \frac{M(\alpha, 0)}{\alpha^{s-1}} \frac{x}{d} \ll \prod_{p < T} \{1 + O(p^{-s/3})\} \frac{x}{d} \ll \frac{x}{d}$$

and so (4.13) can be written in the form

$$(4.32) \quad S_d = \frac{6}{\pi^2} \frac{1}{A_1} \prod_{p|A_1} \left(1 - \frac{1}{p^2}\right)^{-1} \prod_{p|d} \left(1 + \frac{1}{p}\right)^{-1} \left\{ \sum_{\substack{v \in \mathbb{Z}_{A_1} \\ (v, k)=1 \\ v \equiv 0(\alpha)}} f(v) \frac{M(A, v)}{A^{s-1}} \frac{\varphi((A_1, v))}{(A_1, v)} \right\} \frac{x}{d} + \\ + O\left(\frac{x}{d} T^{-s/3} + A^2 x^\delta \sqrt{\frac{x}{d}}\right), \quad \forall \delta > 0.$$



Moreover, since  $(\alpha, k)=1$  and  $(\alpha, A_1)=1$

$$\begin{aligned}
 (4.33) \quad & \sum_{\substack{v \in \mathbb{Z}_{A_1} \\ (v, k)=1, v \equiv 0(\alpha)}} f(v) \frac{M(A, v)}{A^{s-1}} \frac{\varphi((A_1, v))}{(A_1, v)} = \\
 & = \frac{M(\alpha, 0)}{\alpha^{s-1}} \sum_{\substack{u \in \mathbb{Z}_{A_1} \\ (u, k)=1}} f(\alpha u) \frac{M(A_1, \alpha u)}{A_1^{s-1}} \frac{\varphi((A_1, u))}{(A_1, u)} = \\
 & = \frac{M(\alpha, 0)}{\alpha^{s-1}} \sum_{t|A_1/k^2} \frac{\varphi(t)}{t} \sum_{\substack{u \in \mathbb{Z}_{A_1} \\ (u, A_1)=t}} f(\alpha u) \frac{M(A_1, \alpha u)}{A_1^{s-1}}.
 \end{aligned}$$

Now

$$(4.34) \quad \sum_{\substack{u \in \mathbb{Z}_{A_1} \\ (u, A_1)=t}} f(\alpha u) \frac{M(A_1, \alpha u)}{A_1^{s-1}} = \frac{M(t, 0)}{t^{s-1}} \sum_{\substack{q \in \mathbb{Z}_{A_1/t} \\ (q, A_1/t)=1}} f(\alpha t q) \frac{M(A_1/t, \alpha t q)}{(A_1/t)^{s-1}}.$$

Since  $(\alpha t, A_1/t)=1$ , when  $q$  runs in  $\mathbb{Z}_{A_1/t}^*$ , also  $\alpha t q$  runs in  $\mathbb{Z}_{A_1/t}^*$ . Write, for  $\alpha t q \equiv \beta \pmod{A_1/t}$

$$(4.35) \quad \alpha t q = \beta + c A_1/t, \quad \beta \in \mathbb{Z}_{A_1/t}.$$

Now, if  $p|q$ , then  $p|A$ , but  $p \nmid A_1/t$  because  $(q, A_1/t)=1$ ; hence  $p \mid \frac{A}{A_1/t} = \alpha t$  and we have

$$(4.36) \quad f(\alpha t q) = f(\alpha t).$$

By (4.34), (4.35) and (4.36) we get

$$(4.37) \quad \sum_{\substack{u \in \mathbb{Z}_{A_1} \\ (u, A_1)=t}} f(\alpha u) \frac{M(A_1, \alpha u)}{A_1^{s-1}} = f(\alpha t) \frac{M(t, 0)}{t^{s-1}} \sum_{q \in \mathbb{Z}_{A_1/t}^*} \frac{M(A_1/t, q)}{(A_1/t)^{s-1}}$$

and (4.33) becomes

$$\begin{aligned}
 (4.38) \quad & \sum_{\substack{v \in \mathbb{Z}_{A_1} \\ (v, k)=1 \\ v \equiv 0(\alpha)}} f(v) \frac{M(A, v)}{A^{s-1}} \frac{\varphi((A_1, v))}{(A_1, v)} = \\
 & = f(\alpha) \frac{M(\alpha, 0)}{\alpha^{s-1}} \sum_{t|A_1/k^2} \frac{f(t) \varphi(t)}{t} \frac{M(t, 0)}{t^{s-1}} \sum_{q \in \mathbb{Z}_{A_1/t}^*} \frac{M(A_1/t, q)}{(A_1/t)^{s-1}}.
 \end{aligned}$$

We define

$$(4.39) \quad g(t) = \frac{f(t) \varphi(t)}{t} \frac{M(t, 0)}{t^{s-1}}$$

$$(4.40) \quad h(t) = \sum_{q \in \mathbb{Z}_t^*} \frac{M(t, q)}{t^{s-1}}$$

and we can write (4.32) in the form

$$(4.41) \quad S_d = \frac{6}{\pi^2} \frac{1}{A_1} \prod_{p|A_1} \left(1 - \frac{1}{p^2}\right)^{-1} \prod_{p|d} \left(1 + \frac{1}{p}\right)^{-1} \frac{f(\alpha) M(\alpha, 0)}{\alpha^{s-1}} \times \\ \times \left( \sum_{t|A_1/k^2} g(t) h\left(\frac{A_1}{t}\right) \right) \frac{x}{d} + O\left(\frac{x}{d} T^{-s/3} + A^2 x^\delta \sqrt{\frac{x}{d}}\right).$$

The functions  $g(t)$ ,  $h(t)$  are clearly multiplicative (for  $h(t)$ , it suffices to use the chinese remainder's theorem).

Moreover, it is easily seen that, if  $g(t)$  and  $h(t)$  are multiplicative, then also

$$G(d) = \sum_{t|d} \frac{g(t)}{h(t)}$$

is multiplicative.

Then, substituting in (4.41) the value of  $S_1$  obtained by the same formula for  $d=1$ , we get

$$(4.42) \quad S_d = F_T(d) \{S_1 + O(xT^{-s/3} + A^2 x^{1/2+\delta})\} + O\left(\frac{x}{d} T^{-s/3} + A^2 x^\delta \sqrt{\frac{x}{d}}\right)$$

where  $F_T(d)$  is defined by (4.28) and, recalling also the previous remark, it is easily seen that  $F_T(d)$  is multiplicative.

By the definition (4.28), we have

$$(4.43) \quad F_T(p) = \begin{cases} \left(1 - \frac{1}{p^2}\right) \left(1 - \frac{M(p, 0)}{p^s}\right)^{-1} \left(1 + \frac{1}{p}\right)^{-1} \frac{1}{p} \left\{1 + \frac{1}{p} \left(1 - \frac{M(p, 0)}{p^s}\right)^{-1} \times \right. \\ \left. \times \left(1 - \frac{1}{M(p, 0)}\right) \left(1 - \frac{1}{p}\right) \frac{M(p, 0)}{p^{s-1}}\right\}^{-1} & \text{if } p < T \\ \left(1 + \frac{1}{p}\right)^{-1} \frac{1}{p} & \text{if } p > T \end{cases}$$

and, recalling (4.14), we deduce

$$(4.44) \quad F_T(p) = \left(1 + \frac{1}{p}\right)^{-1} \frac{1}{p} + O(p^{-s/3}),$$

from which we obtain

$$(4.45) \quad F_T(d) \ll \frac{1}{d}$$

and also,  $F(d)$  being defined by (2.3)

$$(4.46) \quad |F(d) - F_T(d)| \ll T^{-s/3}.$$

By (4.45), (4.46) formula (4.42) can be also written in the form

$$(4.47) \quad S_d = F(d) S_1 + O\left(\frac{x}{d} T^{-s/3} + A^2 x^\delta \sqrt{\frac{x}{d}}\right).$$

5. In this section, we prove the Theorems stated in Section 2 and we explicitly compute the function  $\delta_x$  of the Bombieri's asymptotic sieve (see (1.11)), that is the distribution function of the primes for our problem. In view of (1.12),  $\delta_x$  determines the distribution function on  $\mathcal{P}_r$  for every  $r$ .

PROOF OF THEOREM 1. We have

$$(5.1) \quad A = k \prod_{p < T} p \simeq ke^T.$$

We choose the parameter  $T$  as follows:

$$(5.2) \quad T = (\log x)^{1/k}.$$

Then, formula (4.47) gives

$$(5.3) \quad S_d = F(d)S_1 + R_d$$

with

$$(5.4) \quad R_d \ll \frac{x}{d} (\log x)^{-s/3k} + x^\delta \sqrt{\frac{x}{d}}$$

for every  $\delta > 0$ . Summing over  $d$ , we get for  $\delta$  small enough

$$(5.5) \quad \sum_{d \leq x^{1-\delta}} |R_d| \ll x (\log x)^{1-(s/3K)}.$$

For  $\delta_x$ , we have, if  $C \rightarrow \infty$  with  $x$

$$(5.6) \quad \sum_{\substack{p < x \\ p \nmid k}} \mathfrak{S}(p) \sim \sum_{p < x} \frac{M(C, p)}{C^{s-1}} \left( 1 - \frac{1}{M(p, 0)} \right) \sim \sum_{p < x} \frac{M(C, p)}{C^{s-1}}$$

where

$$(5.7) \quad C = k \prod_{p < T_1} p \sim ke^{T_1}$$

for  $T_1$  to be chosen in the sequel. Now

$$(5.8) \quad \sum_{p < x} \frac{M(C, p)}{C^{s-1}} = \sum_{\substack{v \in \mathbb{Z}_C \\ (v, C)=1}} \frac{M(C, v)}{C^{s-1}} \sum_{\substack{p < x \\ p \equiv v(C)}} 1.$$

By the classical form of the remainder term of the prime number formula of arithmetic progressions, see for instance [4], we have

$$(5.9) \quad \pi(x, C, v) = \frac{\text{li } x}{\varphi(C)} + O(x \exp(-c \sqrt{\log x}))$$

uniformly for  $C \leq (\log x)^{1/2}$ . So, we obtain by (5.8), (5.9)

$$(5.10) \quad \sum_{p < x} \frac{M(C, p)}{C^{s-1}} = \frac{1}{C^{s-1}} \sum_{v \in \mathbb{Z}_C^*} M(C, v) \left\{ \frac{\text{li } x}{\varphi(C)} + O(x \exp(-c \sqrt{\log x})) \right\}$$

if  $C \equiv (\log x)^{1/2}$ . Now

$$h(C) = \sum_{v \in \mathbb{Z}_C^*} \frac{M(C, v)}{C^{s-1}}$$

is multiplicative; hence

$$(5.11) \quad h(C) = \frac{1}{C^{s-1}} \prod_{p|k} \{p^{2s} - M(p^2, 0)\} \prod_{\substack{p \nmid k \\ p < T_1}} \{p^s - M(p, 0)\}.$$

Recalling (4.14), we have

$$(5.12) \quad h(C) \ll \varphi(C)$$

and so (5.10) becomes

$$(5.13) \quad \sum_{p < x} \frac{M(C, p)}{C^{s-1}} = \frac{h(C)}{\varphi(C)} \operatorname{li} x + O(x \exp(-c \sqrt{\log x}))$$

from which we deduce (see also (5.18) later)

$$(5.14) \quad \sum_{\substack{p < x \\ p \nmid k}} \mathfrak{S}(p) \sim \frac{h(C)}{\varphi(C)} \frac{x}{\log x} \sim R \frac{x}{\log x}$$

with

$$R = \frac{h(k^2)}{\varphi(k^2)} \prod_{p \nmid k} \left(1 - \frac{M(p, 0)}{p^s}\right) \left(1 - \frac{1}{p}\right)^{-1}.$$

(To verify the conditions  $C \rightarrow \infty$ ,  $C < (\log x)^{1/2}$ , it suffices to choose  $T_1 = (\log \log x)^{1/2}$ . By (4.41) with  $d=1$ , we have

$$(5.16) \quad S_1 \sim \frac{6}{\pi^2} \frac{h(A)}{A} \prod_{p|A} \left(1 - \frac{1}{p^2}\right)^{-1} \left\{ \sum_{t|A/k^s} \frac{g(t)}{h(t)} \right\} x$$

and, observing that

$$(5.17) \quad f(p) = 1 + O(p^{-s/3})$$

$$(5.18) \quad g(p) = 1 - \frac{1}{p} + O(p^{-s/3})$$

$$(5.19) \quad h(p) = p \left(1 - \frac{1}{p}\right) + O(p^{-s/3})$$

we obtain

$$(5.20) \quad S_1 = \frac{h(k^2)}{k^2} \prod_{p \nmid k} \frac{h(p)}{p} \left(1 + \frac{g(p)}{h(p)}\right) x + O(x (\log x)^{-s/3k}) = Lx + O(x (\log x)^{-s/3k})$$

with

$$(5.21) \quad L = \frac{h(k^2)}{k^2} \prod_{p \nmid k} \frac{h(p)}{p} \left(1 + \frac{g(p)}{h(p)}\right).$$

Finally, by (1.10), (1.11), (5.14) and (5.19) we obtain

$$(5.22) \quad \delta_x \sim \delta = \frac{k^2}{\varphi(k^2)} \prod_{p|k} \left(1 + \frac{g(p)}{h(p)}\right)^{-1} (1 - F(p))^{-1} = \frac{R}{HL}. \quad \blacksquare$$

PROOF OF THEOREM 2. We define

$$(5.23) \quad \Gamma = \frac{\Gamma^s \left(1 + \frac{1}{k}\right)}{\Gamma\left(\frac{s}{k}\right)}.$$

By (3.23), (3.27) we have

$$(5.24) \quad \sum_{\substack{n < x \\ (n, k) = 1 \\ n \equiv 0(d)}} \mu^2(n) r(n) = \Gamma \sum_{\substack{n < x \\ (n, k) = 1 \\ n \equiv 0(d)}} \mathfrak{S}(n) n^{(s/k)-1} + R'_d$$

with

$$(5.25) \quad R'_d \ll \frac{x^{s/k}}{d} \left(\frac{x}{d}\right)^{-1/k^2}$$

$$(5.26) \quad \sum_{d < x^{1-\varepsilon}} |R'_d| \ll x^{(s/k)-\eta(\varepsilon)}.$$

By partial summation, we have by (5.3), (5.20)

$$(5.27) \quad N_d = \Gamma \sum_{\substack{n < x \\ (n, k) = 1 \\ n \equiv 0(d)}} \mathfrak{S}(n) n^{(s/k)-1} = \Gamma L \frac{k}{s} F(d) x^{s/k} + O\left(\frac{x^{s/k}}{d} (\log x)^{-s/3k}\right).$$

Writing (5.27) for  $d=1$  and substituting in (5.27), we get, recalling also (4.45)

$$(5.28) \quad N_d = N_1 F(d) + O\left(\frac{x^{s/k}}{d} (\log x)^{-s/3k}\right)$$

with

$$(5.29) \quad N_1 = \Gamma L \frac{k}{s} x^{s/k} + O(x^{s/k} (\log x)^{-s/3k})$$

and the estimate for the sum of the error terms (5.28) is clear.

Finally, for the distribution function of the primes, we have by (5.14), using partial summation,

$$(5.30) \quad \sum_{\substack{p < x \\ p \nmid k}} r(p) \sim \Gamma \sum_{p < x} \mathfrak{S}(p) p^{(s/k)-1} \sim \Gamma \frac{k}{s} R \frac{x^{s/k}}{\log x}.$$

By (1.11), (5.29) and (5.30) we have

$$(5.31) \quad \Gamma \frac{k}{s} R = H\Gamma L \frac{k}{s} \delta'$$

from which

$$(5.32) \quad \delta' = \frac{R}{HL} = \delta. \quad \blacksquare$$

#### REFERENCES

- [1] BOMBIERI, E., On the large sieve, *Mathematika* **12** (1965), 201—225. *MR* 33#5590.
- [2] BOMBIERI, E., The asymptotic sieve, *Rend. Accad. Naz.* **40** (1975—76), no. 1/2, 243—269. *MR* 58#10799.
- [3] DAVENPORT, H., *Analytical methods for diophantine equations and diophantine inequalities*, Ann Arbor Publ., Michigan, 1963. *MR* 28#3002.
- [4] DAVENPORT, H., *Multiplicative number theory*, second ed., Graduate Texts in Math., Springer-Verlag, New York—Berlin, 1980. *MR* 82m:10001.
- [5] HALBERSTAM, H. and RICHERT, H. E., *Sieve methods*, Academic Press, London—New York, 1974. *MR* 54#12689.
- [6] IWANIEC, H., A new form of the error term in the linear sieve, *Acta Arith.* **37** (1980), 307—320. *MR* 82d:10069.
- [7] IWANIEC, H., Almost-primes represented by quadratic polynomials, *Invent. Math.* **47** (1978), 171—188. *MR* 58#5553.
- [8] IWANIEC, H., Sieving limits, *Seminar on Number theory*, Paris, 1979—80, 151—165. *MR* 82j:10004.
- [9] SALERNO, S., A note on Selberg sieve (to appear).
- [10] SALERNO, S., Iwaniec's bilinear form of the error term in Selberg sieve (to appear).
- [11] SELBERG, A., Sieve methods, *Proc. Sympos. Pure Math.* **20** (1971), 311—351. *MR* 58#27861.
- [12] VAUGHAN, R. C., *The Hardy—Littlewood method*, Cambridge tracts in math., Cambridge Univ. Press, 1981.
- [13] VINOGRADOV, I. M., *The method of trigonometrical sums in the theory of numbers*, English translation, Interscience, London, 1954. *MR* 15—941.
- [14] WEIL, A., Numbers of solutions of equations in finite fields, *Bull. Amer. Math. Soc.* **55** (1949), 497—508. *MR* 10—592.

( Received December 13, 1983 )

ISTITUTO DI MATEMATICA  
FACOLTA' DI INGEGNERIA  
UNIVERSITA' DI SALERNO  
I—84100 SALERNO  
ITALY

# BEST APPROXIMATION IN $L_p(w)$ BY ALGEBRAIC POLYNOMIALS

JÖRGEN LÖFSTRÖM

## Abstract

We consider the best approximation  $E_p(n, f)$  by algebraic polynomials of degree at most  $n$  on  $L_p(w)$ , where  $w(x) = (1 - x^2)^{\nu - (1/2)}$ ,  $\nu > 0$ ,  $-1 < x < 1$ . We give necessary and sufficient conditions for  $E_p(n, f) = O(n^{-\alpha})$ ,  $n \rightarrow \infty$ , using the theory of orthogonal polynomials and the theory of interpolation spaces.

## 0. Introduction

We shall consider best approximation by algebraic polynomials on  $L_p(w)$  where  $w(x) = (1 - x^2)^{\nu - (1/2)}$ ,  $\nu > 0$ . Thus we shall study the functional

$$(0.1) \quad E_p(n, f) = \inf \left\{ \left( \int_{-1}^1 |f(x) - q(x)|^p w(x) dx \right)^{1/p}, \deg q \leq n \right\}.$$

More precisely we shall characterize the space of all  $f \in L_p(w)$  such that

$$(0.2) \quad E_p(n, f) = O(n^{-\alpha}), \quad n \rightarrow \infty,$$

where  $\alpha$  is a positive real number. This question has been considered by Ky [5] who gave necessary and sufficient conditions for (0.2) in the case  $w=1$  and  $0 < \alpha < 1$ . For interval values of  $\alpha$ , de Vore—Scott proved (for  $p=1$ ) that

$$E_1(n, f) \leq c_\alpha n^{-\alpha} \int_{-1}^1 (1 - x^2)^{\alpha/2} |f^{(\alpha)}(x)| w(x) dx.$$

We shall give necessary and sufficient conditions for (0.2) for all positive real  $\alpha$  and  $1 \leq p \leq \infty$ , in terms of interpolation spaces, (see Theorem 1 in Section 2). Using a generalized modulus of continuity, similar to the one Ky [5] used in the case  $w=1$ ,  $0 < \alpha < 1$ , we make our conditions more explicite. We consider sufficient conditions for (0.2) in Section 3 (see Corollary 1) and in Section 4 (Corollary 2). Converse estimates, giving necessary conditions for (0.2) are given in Section 5. Note that our explicite conditions are both necessary and sufficient if  $\alpha$  is a non-even, positive real number. (The restriction “ $\alpha$  non-even” is removed in the case  $\nu > 1$ .) See Corollary 3, Section 5. Our method of proof relies on the theory of orthogonal polynomials, where Szegös book [9] is our main source. We also use the theory of interpolation spaces as presented in [1].

1980 *Mathematics Subject Classification*. Primary 41A10; Secondary 41A25.

*Key words and phrases*. Best approximation, algebraic polynomials, orthogonal polynomials, interpolation spaces, Gegenbauer or ultraspherical polynomials.



### 1. Preliminaries

Let  $\nu$  be a given positive real number. We shall consider the weight function

$$w(x) = (1-x^2)^{\nu-(1/2)}, \quad -1 < x < 1,$$

and the corresponding differential operator  $A$ , defined by

$$\begin{aligned} (Af)(x) &= -w(x)^{-1} \frac{d}{dx} \left( (1-x^2) w(x) \frac{df}{dx} \right) = \\ &= (2\nu+1)xf'(x) - (1-x^2)f''(x). \end{aligned}$$

Then  $A$  defines a self-adjoint operator on the Hilbert space  $L_2(w)$ . The spectrum of  $A$  is discrete, consisting of the eigenvalues  $\lambda_m = m(m+2\nu)$ , ( $m=0, 1, 2, \dots$ ). The corresponding eigenfunctions are the so called Gegenbauer or ultraspherical polynomials  $P_m$ , given by

$$P_m(x) = c_m w(x)^{-1} \frac{d^m}{dx^m} ((1-x^2)^m w(x)).$$

As a general reference on the polynomials  $P_m$  we use Szegő's book on orthogonal polynomials [9]. Here we shall list a few basic facts needed in what follows.

From the definition we see that  $P_m$  is a polynomial of degree  $m$ . We shall choose the normalization constant  $c_m$  so that

$$(1.1) \quad P_m(1) = 1, \quad (m = 0, 1, 2, \dots).$$

(This means that  $c_m = (-1)^m 2^{-m} \Gamma(\nu+1/2)/\Gamma(m+\nu+1/2)$ .) Then clearly  $P_0(x) = 1$ . Moreover

$$(1.2) \quad \sup_{-1 < x < 1} |P_m(x)| = 1.$$

There is an important recursive formula for  $P_m$ :

$$(1.3) \quad xP_m(x) = A_{m+1}P_{m+1}(x) + B_{m-1}P_{m-1}(x),$$

where

$$A_m = \frac{m+2\nu-1}{2(m+\nu-1)}, \quad m = 1, 2, 3, \dots,$$

$$B_m = \frac{m+1}{2(m+\nu-1)}, \quad m = 0, 1, 2, \dots$$

(If we put  $P_{-1} = 0$ , (1.3) remains valid even in the case  $m=0$ .) Next we consider the eigenfunction expansion

$$f \sim \sum_{m=0}^{\infty} f_m \hat{f}(m) P_m,$$

where (for  $f \in L_1(w)$ ),

$$\hat{f}(m) = \int_{-1}^1 f(x) P_m(x) w(x) dx,$$

and

$$f_m = 1 / \int_{-1}^1 |P_m(x)|^2 w(x) dx = 2^{1-2v} \frac{(m+v)\Gamma(m+2v)}{m! \Gamma\left(v + \frac{1}{2}\right)}.$$

Note that  $r_m = O(m^{2v})$  as  $m \rightarrow \infty$ .

There is an important convolution structure connected with the Gegenbauer polynomials  $P_m$ . In order to define the (generalized) convolution between  $f$  and  $g$  we shall need a (generalized) translation which we shall denote by  $\tau_s$ . Inspired by the formula  $\exp(im(x+s)) = \exp(ims) \cdot \exp(imx)$  we put  $(\tau_s P_m)(x) = P_m(s) P_m(x)$ . Then it is natural to define  $\tau_s f$  by the formula

$$\tau_s f \sim \sum r_m \hat{f}(m) P_m(s) P_m, \quad (-1 < s < 1).$$

By (1.2),  $\tau_s$  is a uniformly bounded operator of  $L_2(w)$ . Note also that  $\tau_1 f = f$  (by (1.1)). It turns out that

$$(1.4) \quad (\tau_s f)(x) = \int_{-1}^1 H(s, x, y) f(y) w(y) dy,$$

where

$$w(s)w(x)w(y)H(s, x, y) = c(1-s^2-x^2-y^2+2sxy)_+^{-1},$$

(with  $c = 2^{1-2v}\Gamma(v)^{-2}$ ). See Bochner [2]. We shall not use the exact expression for the kernel  $H$  but we shall merely use the fact that  $H$  is non-negative, because this implies that  $\tau_s$  is uniformly bounded on  $L_p(w)$ ,  $1 \leq p \leq \infty$ . To see this just note that (1.4) implies

$$\|\tau_s f\|_p \leq \|f\|_p \sup_{-1 < x < 1} \int_{-1}^1 H(s, x, y) w(y) dy$$

and

$$\int_{-1}^1 H(s, x, y) w(y) dy = (\tau_s P_0)(x) = P_0(s) P_0(x) = 1,$$

since  $P_0(x) = 1$ . Thus we have

$$(1.5) \quad \|\tau_s f\|_p \leq \|f\|_p, \quad -1 < s < 1, \quad 1 \leq p \leq \infty.$$

We are now ready to define the (generalized) convolution  $f * g$  between two functions  $f$  and  $g$  in  $L_1(w)$ . The definition is the natural one:

$$(f * g)(x) = \int_{-1}^1 (\tau_s f)(x) g(s) w(s) ds.$$

Then (1.5) implies

$$(1.6) \quad \|f * g\|_p \leq \|f\|_1 \|g\|_p.$$

It is easy to check that

$$(1.7) \quad (f * g)^\wedge(m) = \hat{f}(m) \hat{g}(m).$$

We have now listed all essential facts about the Gegenbauer polynomials needed in the sequel. Let us now conclude this section with some facts on interpolation

spaces which will be essential to us. As a general reference on interpolation spaces we use Bergh—Löfström [1].

If  $X_0$  and  $X_1$  are two semi-normed spaces with seminorms  $\|\cdot\|_0$  and  $\|\cdot\|_1$ , respectively, the interpolation spaces  $(X_0, X_1)_{\theta, \infty}$ ,  $0 < \theta < 1$  can be defined in several equivalent ways. We shall need the following two equivalent definitions, valid in the case when  $X_1$  is continuously embedded in  $X_0$ :

*First definition:* Let  $N$  be any positive integer and suppose  $0 \leq \theta \leq 1$ . Then  $f \in (X_0, X_1)_{\theta, \infty}$  if and only if for every  $t \in ]0, 1[$ , there is a decomposition  $f = f_0 + f_1$  such that

$$\|f_0\|_0 = O(t^{\theta N}), \quad t^N \|f_1\|_1 = O(t^{\theta N}), \quad t \rightarrow 0.$$

*Second definition:* Let  $N$  be any positive integer and suppose  $0 < \theta < 1$ , (strict inequalities). Then  $f \in (X_0, X_1)_{\theta, \infty}$  if and only if there is a decomposition  $f = \sum_{j=0}^{\infty} f_j$  (convergence in  $X_0$ ), such that

$$\|f_1\|_0 = O(2^{-\theta N j}), \quad 2^{-N j} \|f_j\|_1 = O(2^{-\theta N j}), \quad j \rightarrow \infty.$$

These definitions will be used in the case when  $X_0 = L_p = L_p(w)$  and when  $X_1 = D_p(A^N)$ , then domain in  $L_p(w)$  of the operator  $A^N$ , with semi-norm  $\|A^N f\|_p$ . For further details on the definitions of  $(X_0, X_1)_{\theta, \infty}$  see [1] and [3]. Let us just mention two immediate consequences of the so called reiteration theorem:

$$(1.8) \quad ((L_p, D_p(A^N))_{\eta, \infty} D_p(A^N))_{\theta, \infty} = (L_p, D_p(A^N))_{\eta + \theta(1-\eta), \infty}, \quad 0 < \eta, \theta < 1,$$

$$(1.9) \quad (L_p, (L_p, D_p(A^N))_{\eta, \infty})_{\theta, \infty} = (L_p, D_p(A^N))_{\theta \eta, \infty}, \quad 0 < \eta, \theta < 1.$$

## 2. A general convergence result

As mentioned in the introduction we shall study the best approximation  $E_p(n, f)$  in  $L_p(w)$  by algebraic polynomials of degree at most  $n$ . Thus

$$E_p(n, f) = \inf \{ \|f - q\|_p : \deg q \leq n \}.$$

For a given  $f$  we shall construct a sequence  $(q_n)_1$  of polynomials  $q_n$  which are almost best possible. This construction imitates standard constructions in the case of trigonometric approximation.

Let  $\varphi$  be a fixed infinitely differentiable function on the real line such that

$$\varphi(u) = \begin{cases} 1 & \text{for } u \leq 1/2, \\ 0 & \text{for } u \geq 1. \end{cases}$$

We then put

$$\Phi_n = \sum_{m=0}^{\infty} r_m \varphi(m/n) P_m,$$

where  $P_m$  is the Gegenbauer polynomial of degree  $m$  and  $r_m$  is defined in the previous section. Note that the series defining  $\Phi_n$  is finite, since the terms with  $m > n$  vanish. Therefore  $\Phi_n$  is a polynomial of degree at most  $n$ .

For a given  $f \in L_p(w)$  we put

$$q_n = \Phi_n * f = \sum_{m=0}^{\infty} r_m \varphi(m/n) \hat{f}(m) P_m.$$

Clearly,  $q_n$  is a polynomial of degree at most  $n$ . Using the notation of the previous section, we now have the following general convergence result.

**THEOREM 1.** *Let  $N$  be a positive integer and  $\alpha$  real number such that  $0 < \alpha < 2N$ . Then for any  $f \in L_p = L_p(w)$ ,  $(1 \leq p \leq \infty)$ , the following conditions are equivalent:*

- (i)  $E_p(n, f) = O(n^{-\alpha}), \quad n \rightarrow \infty,$
- (ii)  $\|q_n - f\|_p = O(n^{-\alpha}), \quad n \rightarrow \infty,$
- (iii)  $f \in (L_p, D_p(A^N))_{\alpha/(2N), \infty}.$

In the case  $\alpha = 2N$  then (iii) implies (i) and (ii). If  $N = 1$  and  $0 < \alpha < 2$  each one of these conditions is equivalent to

$$(iv) \quad \|\tau_s f - f\|_p = O((1-s)^{\alpha/2}), \quad s \rightarrow 1.$$

The proof of the theorem will be based on the following lemma.

**LEMMA 1.** *Let  $L$  be a non-negative integer and  $\delta$  a positive real number. Assume that  $(\psi_n)_{n=1}^{\infty}$  is a family of infinitely differentiable functions on the real line such that, for some number  $\varrho > 0$ .*

$$|\psi_n(u) - \varrho| \leq c_0 u^{\delta}, \quad 0 < u < 1$$

$$(2.1a) \quad |\psi_n^{(k)}(u)| \leq c_u u^{\delta-k}, \quad 0 < u < 1, \quad k = 1, 2, \dots$$

and

$$(2.1b) \quad |\psi_n^{(k)}(u)| \leq c_k u^{-2L-\delta-k}, \quad u > 1, \quad k = 0, 1, 2, \dots$$

Let  $\Psi_n$  be defined by the relation  $\Psi_n(m) = \psi_n(m/n)$ , i.e.

$$\Psi_n = \sum_{m=0}^{\infty} r_m \psi_n(m/n) P_m.$$

Then

$$\|A^L \Psi_n * f\|_p \leq C n^{2L} \|f\|_p, \quad 1 \leq p \leq \infty,$$

where  $C$  depends on  $\varrho$ ,  $\delta$ ,  $L$  and the constants  $C_k$  appearing in the estimates (2.1).

**PROOF OF THEOREM 1.** The equivalence (iii)  $\Leftrightarrow$  (iv) was proved in Löfström—Peetre [6]. Since we shall not use (iv) in the sequel we do not recall the proof here. Thus we shall just prove (i)  $\Rightarrow$  (ii)  $\Rightarrow$  (iii)  $\Rightarrow$  (i).

If (i) holds there is a sequence  $(e_n)_{n=1}^{\infty}$  of polynomials  $e_n$  such that

$$\|e_n - f\|_p \leq C n^{-\alpha}, \quad \deg e_n \leq n.$$

Since  $\Phi_{2n}(m)=1$  if  $m \leq n$  we have  $\Phi_{2n} * e_n = e_n$ . Consequently,  $q_{2n} - f = \Phi_{2n} * (f - e_n) - (f - e_n)$ . Using Lemma 1 (in the case  $L=0$ ) we conclude that

$$\|q_{2n} - f\|_p \leq C \|f - e_n\|_p \leq Cn^{-\alpha}.$$

This proves the implication (i)  $\Rightarrow$  (ii).

Next assume that (ii) holds. We shall prove (iii) using the second definition given in the previous section. To do that we use an infinitely differentiable function  $\sigma_1$  on the real line with support on  $1/2 \leq u \leq 2$  such that  $\sigma_1(u) > 0$  on  $1/2 < u < 2$ . Then we put  $\sigma(u) = \sigma_1(u) / \sum_{j=-\infty}^{\infty} \sigma(2^{-j}u)$ . Then  $\sigma$  is infinitely differentiable with support on  $1/2 \leq u \leq 2$  and

$$\sum_{j=1}^{\infty} \sigma(m2^{2j}) = 1 \quad \text{if } m = 1, 2, \dots$$

(The construction of  $\sigma$  is a standard construction frequently used in interpolation theory. See for instance [1] ch. 6.) We now put  $L_0 = 1$ , and

$$L_j = \sum_{m=0}^{\infty} r_m \sigma(m2^{-j}) P_m, \quad j = 1, 2, \dots,$$

and

$$f_j = L_j * f.$$

Clearly,  $f_j$  and  $L_j$  are polynomials of degree at most  $2^{j+1}$ . But we also have  $L_j * q_{2^{j-2}} = 0$  since  $L_j(m) = \sigma(m2^{-j}) = 0$  if  $m \leq 2^{j-2}$ . Thus

$$f_j = L_j * (f - q_{2^{j-2}}).$$

Using Lemma 1 we see, however, that

$$(2.2) \quad \|f_j\|_p \leq C \|f - q_{2^{j-2}}\|_p \leq C2^{-\alpha j}.$$

This implies that  $\sum_{j=0}^{\infty} f_j$  converges in  $L_p(w)$ , with the obvious limit  $f$ . Next it is easy to see that

$$A^N(L_j * f) = (A^N L_j) * f,$$

and thus as above

$$A^N f_j = (A^N L_j) * (f - q_{2^{j-2}}).$$

Now Lemma 1 implies that  $\|A^N L_j\|_1 \leq C2^{2Nj}$ . Thus

$$(2.3) \quad 2^{-2Nj} \|A^N f_j\|_p \leq C2^{-\alpha j}.$$

But according to the second definition of Section 1, (2.2) and (2.3) implies  $f \in (L_p, D_p(A^N))_{\alpha/(2N), \infty}$ . We have proved the implication (ii)  $\Rightarrow$  (iii). Assuming that (iii) holds we can find a sequence  $(f_n)_1$  of functions in  $D_p(A^N)$ , such that

$$\|f - f_n\|_p \leq Cn^{-\alpha},$$

$$n^{-2N} \|A^N f_n\|_p \leq Cn^{-\alpha}.$$

This follows from the first definition of Section 1, if we put  $t=1/n$ . Then

$$E_p(n, f) \equiv \|\Phi_n * f - f\|_p \equiv \|\Phi_n * f_n - f_n\|_p + \|f_n - f\|_p.$$

It remains to show that the first term on the right-hand side is bounded by a constant times  $n^{-\alpha}$ . Since the eigenvalues of  $A$  are  $m(m+2v)$ , we have

$$(A^N f_n)^\wedge(m) = (m(m+2v))^N \hat{f}_n(m).$$

Thus

$$\begin{aligned} \Phi_n * f_n - f_n &= \sum_{m=0}^{\infty} r_m (\varphi(m/n) - 1) \hat{f}_n(m) P_m = \\ &= \sum_{m=0}^{\infty} r_m \frac{\varphi(m/n) - 1}{(m(m+2v))^N} (A^N f_n)^\wedge(m) P_m = \\ &= n^{-2N} \sum_{m=0}^{\infty} r_m \Psi_n(m/n) (A^N f_n)^\wedge(m) P_m, \end{aligned}$$

where

$$\psi_n(u) = \frac{\varphi(u) - 1}{(u(u+2v/n))^N}.$$

Now  $\psi_n$  satisfies the assumptions of Lemma 1, (with  $L=0$ ,  $\delta=2N$ ). Thus we conclude that

$$\|\Phi_n * f_n - f_n\|_p \leq C n^{-2N} \|A^N f_n\|_p \leq C n^{-\alpha}.$$

Now we have proved the implication (iii)  $\Rightarrow$  (i).

PROOF OF LEMMA 1. The essential points in the proof of Lemma 1 can be found in Peetre—Vretare [7]. For completeness we prefer, however, to give a direct proof here.

First we note that

$$A^L \Psi_n = n^{2L} \sum r_m \tilde{\psi}_n(m/n) P_m,$$

where

$$\tilde{\psi}_n(u) = (u(u+2v/n))^L \psi_n(u).$$

If (2.1a) and (2.1b) hold for  $\psi_n$  then

$$|\tilde{\psi}_n^{(k)}(u)| \leq C_k u^{\delta-k}, \quad 0 < u < 1, \quad k = 0, 1, 2, \dots,$$

and

$$|\tilde{\psi}_n^{(k)}(u)| \leq \bar{C}_k u^{-\delta-k}, \quad u > 1, \quad k = 0, 1, 2, \dots$$

Therefore it is enough to prove the lemma in the case  $L=0$ . It is no restriction to assume that  $\varrho=0$ , since otherwise we write  $\psi_n(u) = \varphi(u)(\psi_n(u) - \varrho) + (1 - \psi(u))\psi_n(u) + \varrho(\varphi(u) - 1) + \varrho$ , where  $\varphi$  is infinitely differentiable and  $\varphi(u)=1$  on  $u < 1/2$ ,  $\varphi(u)=0$  on  $u > 1$ . The three first terms vanish at  $u=0$  and the fourth is just the constant function for which the result of the lemma (with  $L=0$ ) holds.

With  $\varrho=0$  the assumptions of the lemma can be written

$$(2.4) \quad |\psi_n^{(k)}(u)| \leq C_k \min(u^\delta, u^{-\delta}) u^{-k}, \quad k = 0, 1, 2, \dots$$

We shall prove that

$$\int_{-1}^1 |\Psi_n(s)| w(s) ds \leq C,$$

since this implies the lemma by formula (1.6).

Let  $\sigma$  be the auxiliary function used in the proof of Theorem 1. Then

$$\psi_n(m/n) = \sum_{-\infty}^{\infty} \sigma(m2^{-j}) \psi_n(m/n) = \sum_{-\infty}^{\infty} \chi_j(m).$$

It will be enough to prove that

$$(2.5) \quad \int_{-1}^1 |\Psi_{j,n}(s)| w(s) ds \leq C \min((n2^{-j})^\delta, (n2^{-j})^{-\delta}),$$

if  $\hat{\Psi}_{j,n}(m) = \chi_j(m) = \sigma(m2^{-j}) \psi_n(m/n)$ . In fact, (2.5) implies that

$$\begin{aligned} \int_{-1}^1 |\Psi_n(s)| w(s) ds &\leq \sum_{-\infty}^{\infty} \int_{-1}^1 |\Psi_{j,n}(s)| w(s) ds \leq \\ &\leq C \sum_{-\infty}^{\infty} \min((n2^{-j})^\delta, (n2^{-j})^{-\delta}) \leq C \int_0^{\infty} \min\left(\frac{n}{x}, \frac{x}{n}\right)^\delta \frac{dx}{x} \leq C. \end{aligned}$$

In order to prove (2.5) we shall use the recursive formula (1.3) for the Gegenbauer polynomial, which gives

$$(1-s)\Psi_{j,n}(s) = \sum_{m=0}^{\infty} r_m \chi_j(m)(1-s)P_m(s) = \sum_{m=0}^{\infty} r_m \nabla \chi_j(m) P_m(s),$$

where

$$\nabla \chi_j(m) = \chi_j(m) - \frac{r_{m-1}}{r_m} A_m \chi_j(m-1) - \frac{r_{m+1}}{r_m} B_m \chi_j(m+1).$$

Writing  $A_0 = \chi_j(-1) = 0$  this formula is valid for  $m=0, 1, 2, \dots$ . Using the explicit formulas for  $r_m, A_m, B_m$  given in Section 1, we get

$$(2.6) \quad \nabla_j(m) = -\frac{1}{2} \frac{m+2v}{m+v} \Delta^2 \tau \chi_j(m) - \frac{v}{m+v} \Delta \tau \chi_j(m),$$

where

$$\tau \chi(m) = \chi(m-1), \quad \Delta \chi(m) = \chi(m+1) - \chi(m).$$

Denoting the iterates of  $\nabla$  by  $\nabla^k$  we now have

$$(1-s)^k \Psi_{j,n}(s) = \sum_{m=0}^{\infty} r_m \nabla^k \chi_j(m) P_m.$$



Using Parseval's formula we conclude that

$$(2.7) \quad \int_{-1}^1 ((1-s)^k |\Psi_{j,n}(s)|)^2 w(s) ds \leq C \sum_{m=0}^{\infty} m^{2\nu} |\nabla^k \chi_j(m)|^2,$$

since  $r_m = O(m^{2\nu})$ .

Next we use Cauchy—Schwarz inequality. Let  $\varrho$  be any number between 0 and 1. Then (2.7) implies

$$\begin{aligned} \int_{|1-s| \leq \varrho^2} |\Psi_{j,n}(s)| w(s) ds &\leq C \left( \int_{|1-s| \leq \varrho^2} (1-s)^{-2k} w(s) ds \right)^{1/2} \left( \sum_{m=0}^{\infty} m^{2\nu} |\nabla^k \chi_j(m)|^2 \right)^{1/2} \leq \\ &\leq C \varrho^{-2k+\nu+(1/2)} \left( \sum_{m=0}^{\infty} m^{2\nu} |\nabla^k \chi_j(m)|^2 \right)^{1/2}. \end{aligned}$$

Here we assume  $k > \nu + 1/2$ . Similarly

$$\int_{|1-s| \leq \varrho^2} |\Psi_{j,n}(s)| w(s) ds \leq C \varrho^{\nu+(1/2)} \left( \sum_{m=0}^{\infty} m^{2\nu} |\chi_j(m)|^2 \right)^{1/2}.$$

We now choose  $\varrho$  so that

$$\varrho^{-2k} \left( \sum_{m=0}^{\infty} m^{2\nu} |\nabla^k \chi_j(m)|^2 \right)^{1/2} \sim \left( \sum_{m=0}^{\infty} m^{2\nu} |\chi_j(m)|^2 \right)^{1/2}.$$

Then we get

$$\begin{aligned} (2.8) \quad &\int_{-1}^1 |\Psi_{j,n}(s)| w(s) ds \leq \\ &\leq C \left( \sum_{m=0}^{\infty} m^{2\nu} |\chi_j(m)|^2 \right)^{(1-\theta)/2} \left( \sum_{m=0}^{\infty} m^{2\nu} |\nabla^k \chi_j(m)|^2 \right)^{\theta/2}, \end{aligned}$$

where

$$\theta = \frac{\nu + \frac{1}{2}}{2k}.$$

It remains to prove

$$\begin{aligned} (2.9) \quad &\left( \sum_{m=0}^{\infty} m^{2\nu} |\nabla^k \chi_j(m)|^2 \right)^{1/2} \leq \\ &\leq C_R 2^{(\nu+(1/2)-2k)j} \min((n2^{-j})^\delta, (n2^{-j})^{-\delta}), \quad k = 0, 1, 2, \dots \end{aligned}$$

In fact, this estimate combined with (2.8) implies (2.5) since

$$\begin{aligned} &\int_{-1}^1 |\Psi_{j,n}(s)| w(s) ds \leq \\ &\leq C 2^{(\nu+(1/2))j} \min((n2^{-j})^\delta, (n2^{-j})^{-\delta}) 2^{-2kj\theta} = C \min((n2^{-j})^\delta, (n2^{-j})^{-\delta}). \end{aligned}$$

We carry out the details in the proof of (2.9) in the case  $k=1$ . By formula (2.6) we see that the left-hand side of (2.9) can be estimated by a constant times the sum of the integrals

$$I_i = \left( \int_0^\infty u^{2v} |u^{-i} \chi_j^{(2-i)}(u)|^2 du \right)^{1/2}, \quad i = 0, 1.$$

Now, by (2.4), we have

$$|\chi^{(i)}(u)| \leq C 2^{-j i} \min((n 2^{-j})^\delta, (n 2^{-j})^{-\delta}).$$

Thus

$$\begin{aligned} I_i &\leq C \left( \int_{2^{j-1}}^{2^{j+1}} u^{2v} du \right)^{1/2} 2^{-2j} \min((n 2^{-j})^\delta, (n 2^{-j})^{-\delta}), \\ &\leq C 2^{(v+(1/2)-2)j} \min((n 2^{-j})^\delta, (n 2^{-j})^{-\delta}). \end{aligned}$$

For the general case  $k > 0$  the argument is quite similar starting with the estimate

$$|\nabla^k \chi_j(m)| \leq C_k \sum_{r=1}^{2k} m^{r-2k} |\Delta^r \tau^k \chi_j(m)|$$

(cf. Peetre—Vretare [7]). We leave the details to the reader. The case  $k=0$  follows at once from the estimate

$$m^{2v} |\chi(m)|^2 \leq C \left( \int_m^{m+1} u^{2v} |\chi(u)|^2 du + \int_m^{m+1} u^{2v} |\chi'(u)|^2 du \right).$$

This completes the proof of Lemma 1.

### 3. Explicite convergence results

In Theorem 1 we gave a complete characterization of the space of all  $f \in L_p = L_p(w)$  such that  $E_p(n, f) = O(n^{-\alpha})$ , in terms of generalized convolutions and interpolation spaces. However, these conditions are not very explicite. The reader will probably have some difficulties in checking the conditions (ii)—(iv) of Theorem 1 even for rather simple functions as  $(1-x^2)^\beta$ . In this section we shall therefore seek for more explicite conditions.

Let us denote by  $\mathcal{P}_k$  the space of all polynomials of degree at most  $k$ . If we write “mod  $\mathcal{P}_k$ ” in a formula we mean that the formula holds if we subtract a suitable polynomial from the function involved.

We shall work with the operator  $B$  defined by

$$Bf(x) = i \sqrt{1-x^2} f'(x).$$

The nice thing about this operator is that

$$(3.1) \quad (Bf)^* = D(f^*)$$

where

$$f^*(\theta) = f(\cos \theta)$$

and  $D=id/d\theta$  denotes ordinary differentiation with respect to  $\theta$ . We shall use the notation  $D_p(B^L)$  for the domain of the operator  $B^L$ , semi-normed by  $\|B^L f\|_p$ .

We shall need the following lemma:

LEMMA 2. Let  $p=1$ . Then

$$\|Af\|_1 \leq C \|B^2 f\|_1 \pmod{\mathcal{P}_1}.$$

PROOF. Without loss of generality we can assume that  $f(0)=0$  and  $f'(\pi/2)=0$ , or equivalently  $f^*(\pi/2)=Df^*(\pi/2)=0$ . Now

$$Af(x) - B^2 f(x) = 2vx f'(x) = \frac{2vx}{i\sqrt{1-x^2}} Bf(x).$$

Therefore it is enough to prove

$$\int_{-1}^1 \frac{1}{\sqrt{1-x^2}} |Bf(x)| w(x) dx \leq C \|B^2 f\|_1.$$

Changing variables ( $x=\cos \theta$ ) this is a consequence of

$$(3.2) \quad \int_{\pi/2}^{\pi} (\sin \theta)^{-1} |Df^*(\theta)| (\sin \theta)^{2\nu} d\theta \leq C \int_{\pi/2}^{\pi} |D^2 f^*(\theta)| (\sin \theta)^{2\nu} d\theta,$$

and from the corresponding inequality involving integrals over the interval  $]0, \pi/2[$ . Now (3.2) is easy to prove by means of a direct calculation which runs like this:

$$\begin{aligned} \int_{\pi/2}^{\pi} (\sin \theta)^{-1} |Df^*(\theta)| (\sin \theta)^{2\nu} d\theta &\leq \int_{\pi/2}^{\pi} (\sin \theta)^{-1} \int_{\pi/2}^{\theta} |D^2 f^*(\eta)| d\eta (\sin \theta)^{2\nu} d\theta = \\ &= \int_{\pi/2}^{\pi} \int_{\eta}^{\pi} (\sin \theta)^{2\nu-1} d\theta |D^2 f^*(\eta)| d\eta \leq C \int_{\pi/2}^{\pi} |D^2 f^*(\eta)| (\sin \eta)^{2\nu} d\eta. \end{aligned}$$

This completes the proof.

We shall now consider a family of operators  $T(s)$ , defined for  $0 < s < \pi/4$  by the formula

$$(T(s)f^*)(\theta) = \begin{cases} f^*(\theta+s) & \text{if } 0 < \theta < \pi/2, \\ f^*(\theta-s) & \text{if } \pi/2 < \theta < \pi. \end{cases}$$

Let  $L_p^*$  denote the space normed by

$$\|f^*\|_{L_p^*} = \left( \int_0^{\pi} |f^*(\theta)|^p (\sin \theta)^{2\nu} d\theta \right)^{1/p}.$$

Then

$$(3.3) \quad \|f\|_p = \|f^*\|_{L_p^*}.$$

Moreover we have

$$(3.4) \quad \|T(s)f^*\|_{L_p^*} \leq 2^\nu \|f^*\|_{L_p^*}.$$

This follows at once from the estimate

$$\sin \theta \cong \begin{cases} \sqrt{2} \sin(\theta+s) & \text{if } 0 < \theta < \pi/2, 0 < s < \pi/4, \\ \sqrt{2} \sin(\theta-s) & \text{if } \pi/2 < \theta < \pi, 0 < s < \pi/4. \end{cases}$$

For a given  $f \in L_p$  we shall now consider the moduli of continuity

$$\omega_p^k(t, f) = \sup_{0 < s < t} \left( \int_0^\pi |A^k(s) f^*(\theta)|^p (\sin \theta)^{2\nu} d\theta \right)^{1/p},$$

$$A^k(s) = \sum_{j=0}^k \binom{k}{j} (-1)^{k-j} T(sj),$$

where  $0 < t < \pi/4k$ . More explicitly we have for instance

$$\begin{aligned} \omega_p^1(t, f) = & \sup_{0 < s < t} \left( \int_0^{\pi/2} |f^*(\theta+s) - f^*(\theta)|^p (\sin \theta)^{2\nu} d\theta + \right. \\ & \left. + \int_{\pi/2}^\pi |f^*(\theta-s) - f^*(\theta)|^p (\sin \theta)^{2\nu} d\theta \right)^{1/p}, \end{aligned}$$

and

$$\begin{aligned} \omega_p^2(t, f) = & \sup_{0 < s < t} \left( \int_0^{\pi/2} |f^*(\theta+2s) - 2f^*(\theta+s) + f^*(\theta)|^p (\sin \theta)^{2\nu} d\theta + \right. \\ & \left. + \int_{\pi/2}^\pi |f^*(\theta-2s) - 2f^*(\theta-s) + f^*(\theta)|^p (\sin \theta)^{2\nu} d\theta \right)^{1/p}. \end{aligned}$$

**LEMMA 3.** Let  $k$  be a positive integer and assume that  $0 < \alpha \leq k$ . Then for  $1 \leq p \leq \infty$ , the following conditions are equivalent

- (i)  $f \in (L_p, D_p(B^k))_{\alpha/k, \infty}$
- (ii)  $\omega_p^k(t, f) = O(t^\alpha), \quad t \rightarrow 0.$

**PROOF.** Using (3.1) and (3.3) it is easily seen to be enough to prove

$$(3.5) \quad f^* \in (L_p^*, D_p^*(D^k))_{\alpha/k, \infty} \Leftrightarrow \sup_{0 < s < t} \|A^k(s) f^*\|_{L_p^*} = O(t^\alpha).$$

We shall prove (3.5) in the case  $p \neq \infty$ . When  $p = \infty$  the result is well-known. (See for instance [3], Theorem 3.4.3, which also gives a model for the following proof.)

We shall start by proving the implication from the right-hand side to the left-hand side of (3.5) in the case  $k=1$ . Consider the isometric operator on  $L_p^*$  defined by

$$(Jf^*)(\theta) = \begin{cases} f^*(\theta) & \text{if } 0 < \theta < \pi/2, \\ -f^*(\theta) & \text{if } \pi/2 < \theta < \pi. \end{cases}$$

Then  $J^2 = I$  and

$$\frac{d}{du} T(u) f^* = JDT(u) f^* = JT(u) Df^*, \quad (D = d/d\theta).$$

For small values of  $t$  (say  $0 < t < \pi/4$ ) we put

$$f_1^* = \frac{1}{t} \int_0^t T(u) f^* du \quad \text{and} \quad f_0 = -\frac{1}{t} \int_0^t \Delta^1(u) f^* du.$$

Then  $f^* = f_0^* + f_1^*$  and

$$\|f_0^*\|_{L_p^*} \leq \sup_{0 < s < t} \|\Delta^1(s) f^*\|_{L_p^*}.$$

Moreover,

$$tJDf_1^* = \int_0^t \frac{d}{du} T(u) f^* du = \Delta^1(t) f^*,$$

and therefore

$$t\|Df_1^*\|_{L_p^*} \leq \sup_{0 < s < t} \|\Delta^1(t) f^*\|_{L_p^*}.$$

This proves the implication in the left direction of (3.5) when  $k=1$ . In the general case ( $k \geq 1$ ) we put for small  $t$

$$f_1^* = -\frac{1}{t^k} \int_0^t \dots \int_0^t \sum_{j=1}^k \binom{k}{j} (-1)^j T((u_1 + \dots + u_k)j/k) f^* du_1 \dots du_k,$$

and

$$f_0^* = \frac{(-1)^k}{t^k} \int_0^t \dots \int_0^t \Delta^k((u_1 + \dots + u_k)/k) f^* du_1 \dots du_k.$$

Then we have again that  $f^* = f_0^* + f_1^*$  and

$$\|f_0^*\|_{L_p^*} \leq \sup_{0 < s < t} \|\Delta^k(s) f^*\|_{L_p^*}.$$

Moreover we have that

$$\begin{aligned} t^k J^k D^k f_1^* &= -\sum_{j=1}^k \binom{k}{j} (-1)^j \left(\frac{k}{j}\right)^k \int_0^t \dots \int_0^t \frac{\partial}{\partial u_1} \dots \frac{\partial}{\partial u_k} T((u_1 + \dots + u_k)j/k) f^* du_1 \dots du_k = \\ &= -\sum_{j=1}^k \binom{k}{j} (-1)^j \left(\frac{k}{j}\right)^k \Delta^k(tj/k) f^*, \end{aligned}$$

and hence

$$t^k \|D^k f_1^*\|_{L_p^*} \leq \sup_{0 < s < t} \|\Delta^k(s) f^*\|_{L_p^*}.$$

This concludes the proof of the implication in the left direction of (3.5).

The converse implication is simpler. Indeed, if  $f^* = f_0^* + f_1^*$  then we get, using (3.4),

$$\|\Delta^k(s) f^*\|_{L_p^*} \leq C \|f_0^*\|_{L_p^*} + \|\Delta^k(s) f_1^*\|_{L_p^*}.$$

But now we have

$$\begin{aligned} \Delta^k(s)f_1^* &= \int_0^s \dots \int_0^s \frac{\partial}{\partial u_1} \dots \frac{\partial}{\partial u_k} T(u_1 + \dots + u_k) f_1^* du_1 \dots du_k = \\ &= \int_0^s \dots \int_0^s J^k T(u_1 + \dots + u_k) D^k f_1^* du_1 \dots du_k. \end{aligned}$$

Thus

$$\|\Delta^k(s)f_1^*\|_{L_p^*} \leq C s^k \|D^k f_1^*\|_{L_p^*}.$$

Using the definition of the space  $(L_p^*, D_p^*(D^k))_{\alpha/k, \infty}$  we now get the implication in the right direction of (3.5). The proof of Lemma 3 is complete.

**THEOREM 2.** Assume that  $k=1$  or  $k=2$ . Let  $0 < \alpha \leq k$  and  $1 \leq p \leq \infty$ . Then

$$\omega_p^k(t, f) = O(t^\alpha), \quad t \rightarrow 0,$$

implies

$$E_p(n, f) = O(n^{-\alpha}), \quad n \rightarrow \infty.$$

**REMARK.** In the case  $w=1$  and  $k=1$  this was proved by Ky [5].

**PROOF.** By Lemma 2 we have (modulo polynomials)

$$D_1(B^2) \subset D_1(A).$$

We shall now see that

$$(3.6) \quad D_1(B) \subset (L_1, D_1(B^2))_{1/2, \infty},$$

or equivalently

$$D_1^*(D) \subset (L_1^*, D_1^*(D^2))_{1/2, \infty}.$$

This is easy to prove, using the construction of the proof of Lemma 3. Indeed, define  $f_1$  and  $f_0$  as in the first part of that proof with  $k=2$ . Then easily

$$t^2 \|D^2 f_1^*\|_{L_1^*} \leq Ct \|Df^*\|_{L_1^*}, \quad \|f_0^*\|_{L_1^*} \leq Ct \|Df^*\|_{L_1^*},$$

which gives the result.

Using Lemma 2 and (3.6) we get

$$D_1(B) \subset (L_1, D_1(A))_{1/2, \infty}$$

and (obviously)

$$D_1(B^2) \subset (L_1, D_1(A))_{1, \infty}.$$

Thus Theorem 1 implies

$$(3.7) \quad \|q_n - f\|_1 \leq Cn^{-1} \|Bf\|_1$$

and

$$(3.8) \quad \|q_n - f\|_1 \leq Cn^{-2} \|B^2 f\|_1.$$

Next we consider the case  $p=\infty$ . Then the  $L_\infty$ -norm is independent of  $v$ . We may therefore take  $v=0$ . Making the standard change of variables  $x=\cos \theta$  the desired result will then reduce to a well-known result on best uniform (non-weighted) approximation by trigonometric polynomials. Indeed we then get

$$\omega_\infty^k(t, f) = O(t^\alpha) \Rightarrow E_\infty(n, f) = O(n^{-\alpha})$$

for all  $k$  and  $0 < \alpha \leq k$ . There is even a converse implication in the case  $0 < \alpha < k$ . But now

$$\omega_{\infty}^1(t, f) \leq Ct \|Bf\|_{\infty},$$

and

$$\omega_{\infty}^2(t, f) \leq Ct^2 \|B^2 f\|_{\infty}.$$

Thus Theorem 1 implies

$$(3.9) \quad \|q_n - f\|_{\infty} \leq Cn^{-1} \|Bf\|_{\infty}$$

and

$$(3.10) \quad \|q_n - f\|_{\infty} \leq Cn^{-2} \|B^2 f\|_{\infty}.$$

Now we use the Riesz—Thorin interpolation theorem (Theorem 5.1.1 in [1]) on (3.7), (3.9) and (3.8), (3.10) to get

$$\|q_n - f\|_p \leq Cn^{-1} \|Bf\|_p,$$

$$\|q_n - f\|_p \leq Cn^{-2} \|B^2 f\|_p.$$

Theorem 1 then implies

$$D_p(B) \subset (L_p, D_p(A))_{1/2, \infty}$$

and

$$D_p(B^2) \subset (L_p, D_p(A))_{1, \infty}.$$

Using (1.6) we now conclude (for  $k=1, 2$ ) that

$$(L_p, D_p(B^k))_{\alpha/k, \infty} \subset (L_p, (L_p, D_p(A))_{k/2, \infty})_{\alpha/k, \infty} = (L_p, D_p(A))_{\alpha/2, \infty}$$

from which the desired result follows by Lemma 3 and Theorem 1.

**COROLLARY 1.** Suppose that  $k=1, 2$  and  $0 < \alpha \leq k$ . Let  $L$  be a non-negative integer. Then

$$\omega_p^k(t, A^L f) = O(t^{\alpha}) \Rightarrow E_p(n, f) = O(n^{-2L-\alpha}).$$

**PROOF.** If  $\omega_p^k(t, A^L f) = O(t^{\alpha})$  then

$$A^L f \in (L_p, D_p(A))_{\alpha/2, \infty},$$

by the proof of Theorem 2. We now use the following

**LEMMA 4.** Suppose that  $0 < \theta \leq 1$ . Let  $L$  and  $N$  be non-negative integers. Then

$$A^L f \in (L_p, D_p(A)^N)_{\theta, \infty} \Rightarrow f \in (L_p, D_p(A^{L+N}))_{\eta, \infty}$$

where

$$\eta = \frac{\theta N + L}{N + L}.$$

With  $N=1$ ,  $\theta=\alpha/2$  this gives, using Theorem 1,

$$A^L f \in (L_p, D_p(A))_{\alpha/2, \infty} \Rightarrow f \in (L_p, D_p(A^{L+1}))_{(\alpha+2L)/(2(1-L)), \infty} \Rightarrow E_p(n, f) = O(n^{-2L-\alpha}).$$

It remains to prove Lemma 4. In order to do that we use Theorem 1. Put  $q_n = \Phi_n * f$  and  $\tilde{q}_n = \Phi_n * A^L f$ . Then

$$A^L f \in (L_p, D_p(A^N))_{\theta, \infty} \Rightarrow \|\tilde{q}_n - A^L f\|_p = O(n^{-2N\theta}).$$



But now

$$q_n - f = n^{-2L} \sum r_m q_n(m/n) (\varphi(m/n) - 1) (A^L f)^\wedge(m) P_m$$

where

$$q_n(u) = \begin{cases} (u(u+2/n))^{-L} & \text{if } u \geq 1/2 \\ 0 & \text{if } u < 1/4. \end{cases}$$

We can choose  $q_n$  so that it is infinitely differentiable. Then Lemma 1 implies

$$\|q_n - f\|_p \leq C n^{-2L} \|\tilde{q}_n - A^L f\|_p \leq C n^{-2(N\theta+L)},$$

which implies  $f \in (L_p, D_p(A^{N+L}))_{\theta, \infty}$ , by Theorem 1. This completes the proof of Lemma 4.

We conclude this section with a discussion on the possibility of extending Theorem 2 to  $k=3, 4, \dots$ . Suppose we wanted to use our method of proof to the case  $k=4$ . Then we would have to prove the implication (in the case  $p=1$ )

$$f \in (L_1, D_1(B^4))_{\alpha/4, \infty} \Rightarrow f \in (L_1, D_1(A^2))_{\alpha/4, \infty}.$$

However, this would call for the inclusion  $D_1(B^4) \subset D_1(A^2)$ , which is *not* true for all values of  $v$ . However, we can prove that if  $v > N-1$  then

$$(3.11) \quad \|A^N f\|_1 \leq C_N \|B^{2N} f\|_1 \pmod{P_{2N-1}}.$$

Using this estimate and the same method of proof as we used in Theorem 2, we can show the following result.

**THEOREM 2'.** Assume that  $v > N-1$ ,  $k=1, 2, \dots, 2N$  and  $0 < \alpha \leq k$ . Then, for  $L=0, 1, 2, \dots$ ,

$$\omega_p^k(t, A^L f) = O(t^\alpha),$$

implies

$$E_p(n, f) = O(n^{-\alpha-2L}).$$

We leave the details of the proof of Theorem 2' to the reader. Instead let us look closer at the proof of (3.11). First we take  $N=2$ . Since

$$(A - B^2)f(x) = \frac{2vx}{i\sqrt{1-x^2}} Bf(x)$$

we have

$$\begin{aligned} (A^2 - B^4)f(x) &= \frac{2vx}{i\sqrt{1-x^2}} B(A + B^2)f(x) = \\ &= \frac{2vx}{i\sqrt{1-x^2}} B\left(\frac{2x}{i\sqrt{1-x^2}} B + 2B^2\right)f(x) = \\ &= \frac{2vx}{i\sqrt{1-x^2}} \left(\frac{2v}{1-x^2} B + \frac{2vx}{i\sqrt{1-x^2}} B^2 + 2B^3\right)f(x). \end{aligned}$$

Now (3.7) for  $N=1$  implies

$$\begin{aligned} \int_0^\pi \left| \frac{\cos \theta}{\sin \theta} \right| \left| \frac{2\nu}{\sin^2 \theta} Df^*(\theta) + \frac{2 \cos \theta}{\sin \theta} D^2 f^*(\theta) + 2D^3 f^*(\theta) \right| (\sin \theta)^2 d\theta &\leq \\ &\leq C \int_0^\pi |D^4 f^*(\theta)| (\sin \theta)^2 d\theta. \end{aligned}$$

Taking  $Df^*(\theta) = \cos \theta$  we see that this implies

$$\int_0^\pi (\sin \theta)^{2\nu-3} d\theta < \infty,$$

i.e.  $\nu > 1$ . Thus (3.11) is not valid for all  $\nu$ . In order to prove (3.7) we use induction to prove

$$(A^N - B^{2N})f(x) = \sum_{\substack{1 \leq k \leq 2N-1 \\ 2 \leq k+j \leq 2N}} c_{i,j,k} x^i (1-x^2)^{-j/2} B^k f(x),$$

We leave the details to the reader. Then (3.11) will follow if we can show

$$(3.12) \quad \int_{-1}^1 ((1-x^2)^{-j/2} |B^k f(x)| w(x) dx \leq C_{j,k} \|B^L f\|_1 \pmod{P_{L-1}}$$

if  $2 \leq k+j \leq L$ ,  $1 \leq k \leq L-1$  and if  $2\nu-j > -1$ . (Note that if  $L=2N$  then  $\nu > N-1$  implies  $2\nu-j > -1$  for all  $j$  in question.) In order to show (3.12) we assume that  $f(0) = \dots = f^{(L-1)}(0) = 0$ . Then (3.12) will follow if we can show that

$$(3.13) \quad \int_{\pi/2}^\pi |D^k f^*(\theta)| (\sin \theta)^{2\nu-j} d\theta < C \int_{\pi/2}^\pi |D^L f^*(\theta)| (\sin \theta)^{2\nu} d\theta,$$

and a similar estimate for the integral over  $]0, \pi/2[$ . Now the left-hand side can be estimated by a constant times

$$\int_{\pi/2}^\pi \int_{\pi/2}^\theta (\theta-\eta)^{L-k-1} |D^L f^*(\eta)| d\eta (\sin \theta)^{2\nu-j} d\theta.$$

By a direct computation this will easily be estimated by the right-hand side of (3.9).

#### 4. Some additional explicite results

In this section we shall consider the operators

$$U_L f(x) = (i\sqrt{1-x^2})^L f^{(L)}(x).$$

LEMMA 5. For  $1 \leq p \leq \infty$ ,  $N=1, 2, \dots$  we have

$$\begin{aligned} \|A^N f\|_p &\leq C_N \|U_{2N} f\|_p \pmod{\mathcal{P}_{2N-1}} \\ \|BA^{N-1} f\|_p &\leq C_N \|U_{2N-1} f\|_p \pmod{\mathcal{P}_{2N-2}} \\ \|A^N f\|_p &\leq C_N \|AU_{2N-2} f\|_p \pmod{\mathcal{P}_{2N-1}}. \end{aligned}$$

PROOF. We start with the formula

$$(4.1) \quad A^N f(x) - U_{2N} f(x) = \sum_{\substack{1 \leq k \leq 2N-1 \\ k-j \leq N}} c_{i,j,k} x^i (1-x^2)^j f^{(k)}(x),$$

which is easy to prove by induction. From this we see that the result follows if we can prove

$$(4.2) \quad \left( \int_{-1}^1 ((1-x^2)^j |f^{(k)}(x)|)^p w(x) dx \right)^{1/p} \leq C \|U_L f\|_p \pmod{\mathcal{P}_{L-1}}$$

if  $k-j \leq L/2$ ,  $1 \leq k \leq L-1$ . By the Riesz—Thorin interpolation theorem we see that it is enough to consider the cases  $p=1$  and  $p=\infty$ . We can assume that  $f(0)=f'(0)=\dots=f^{(L-1)}(0)=0$ . Then

$$\begin{aligned} \int_0^1 (1-x^2)^j |f^{(k)}(x)| w(x) dx &\leq C \int_0^1 \int_y^1 (1-x)^{j+v-(1/2)} (x-y)^{L-k-1} dx |f^{(L)}(y)| dy \leq \\ &\leq C \int_0^1 (1-y)^{L+j-k+v-(1/2)} |f^{(L)}(y)| dy. \end{aligned}$$

The right-hand side is bounded by a constant times  $\|U_L f\|_1$  since  $L+j-k \geq 0$ . This settles the case  $p=1$ . If  $p=\infty$  we just note that

$$\begin{aligned} (1-x^2)^j |f^{(k)}(x)| &\leq \|U_L f\|_\infty (1-x)^j \int_0^x (x-y)^{L-k-1} (1-y)^{-L/2} dy \leq \\ &\leq C \|U_L f\|_\infty \quad \text{if } 0 < x < 1. \end{aligned}$$

This proves the first part of Lemma 6. The remaining parts are proved in the same way. Using (4.1) it is easy to see that

$$BA^{N-1} f(x) - U_{2N-1} f(x) = \sum_{\substack{1 \leq k \leq 2N-2 \\ k-j \leq N-1}} c_{i,j,k} x^i (1-x^2)^j f^{(k)}(x).$$

From (4.2) we now get the second conclusion of the lemma. The third one follows from an analogous formula for  $A(A^{N-1} - U_{2N-2})$  which the reader will have no difficulty to derive.

COROLLARY 2. For  $1 \leq p \leq \infty$ ,  $L=1, 2, \dots$ , we have

$$U_L f \in L_p \Rightarrow E_p(n, f) = O(n^{-L}).$$

For  $0 < \alpha \leq k$ ,  $k=1, 2$  and  $L$  even we also have

$$\omega_p^k(t, U_L f) = O(t^{-\alpha}) \Rightarrow E_p(n, f) = O(n^{-L-\alpha}).$$

REMARK. The first result was proved in de Vore—Scott [4].

PROOF. If  $L=2N$  the first result follows at once from Lemma 5 and Theorem 1. If  $L=2N-1$ , Lemma 6 gives  $U_{2N-1} f \in L_p \Rightarrow BA^{N-1} f \in L_p$ . Hence  $\omega_p^1(t, A^{N-1} f) = O(t)$ . Now Corollary 1 gives the desired result.

In order to prove the second part with  $L=2N-2$  we use Lemma 2.5 and Corollary 1 to get

$$\begin{aligned} U_{2N-2}f \in L_p &\Rightarrow A^{N-1}f \in L_p \\ BU_{2N-2}f \in L_p &\Rightarrow BA^{N-1}f \in L_p \Rightarrow A^{N-1}f \in (L_p, D_p(A))_{1/2, \infty} \\ B^2U_{2N-2}f \in L_p &\Rightarrow AU_{2N-2}f \in L_p \Rightarrow A^{N-1}f \in D_p(A). \end{aligned}$$

Thus

$$U_{2N-2}f \in (L_p, D_p(B^k))_{\alpha/k, \infty} \Rightarrow A^{N-1}f \in (L_p, D_p(A))_{\alpha/2, \infty}.$$

Using Lemma 3, 4 and Theorem 1 the result now follows.

## 5. Converse results

In this section we shall use two versions of Markov's inequality to derive results which are converses of the results given in Theorem 2 and Corollary 1.

LEMMA 6 (Markov inequality). *Suppose that  $1 \leq p \leq \infty$ . Let  $q$  be a polynomial of degree at most  $n$ . Then*

$$\|B^k q\|_p \leq C_k n^k \|q\|_p,$$

and

$$\|A^L q\|_p \leq C_L n^{2L} \|q\|_p.$$

This lemma is proved in Stein [8]. (The second result will also follow easily from Lemma 1.)

THEOREM 3. *We have the following converse implications valid for  $L=0, 1, 2, \dots$  and  $1 \leq p \leq \infty$ :*

$$0 < \alpha < 1: E_p(n, f) = O(n^{-\alpha-2L}) \Rightarrow \omega_p^1(t, A^L f) = O(t^\alpha),$$

$$0 < \alpha < 2: E_p(n, f) = O(n^{-\alpha-2L}) \Rightarrow \omega_p^2(t, A^L f) = O(t^\alpha).$$

More generally if  $k=1, 2, 3, \dots$  we have

$$0 < \alpha < k: E_p(n, f) = O(n^{-\alpha-2L}) \Rightarrow \omega_p^k(t, A^L f) = O(t^\alpha).$$

PROOF. If  $E_p(n, f) = O(n^{-\alpha-2L})$  there is a sequence  $(e_n)_1^\infty$  of polynomials of degree at most  $n$  so that  $\|f - e_n\|_p \leq C n^{-\alpha-2L}$ . Then  $f = \sum_{j=0}^\infty f_j$  if  $f_j = e_{2^j+1} - e_{2^j}$  for  $j=1, 2, 3, \dots$  and if  $f_0 = e_1$ . Since

$$\|f_j\|_p \leq C 2^{-(\alpha+2L)j},$$

we get from Lemma 6 that

$$\|A^L f_j\|_p \leq C 2^{2Lj} \|f_j\|_p \leq C 2^{-\alpha j}$$

and

$$2^{kj} \|B^k A^L f_j\|_p \leq C \|A^L f_j\|_p \leq C 2^{-\alpha j}.$$

Using the second definition of interpolation spaces we get

$$A^L f \in (L_p, D_p(B^k))_{\alpha/k, \infty}$$

if  $0 < \alpha < k$ . Hence Lemma 3 implies

$$\omega_p^k(t, A^L f) = O(t^\alpha).$$

This proves the theorem.

The following corollary summarizes the contents of Corollary 1, and Theorems 2' and 3.

**COROLLARY 3.** *Suppose that  $1 \leq p \leq \infty$  and  $\beta > 0$ .*

(i) *If  $\nu > 0$  and  $\beta$  is non-even then*

$$E_p(n, f) = O(n^{-\beta}) \Leftrightarrow \omega_p^2(t, A^L f) = O(t^{\beta-2L})$$

*provided that we choose the integer  $L$  so that*

$$2L < \beta < 2L + 2.$$

(ii) *If  $\nu > 1$  and  $\beta$  arbitrary, then*

$$E_p(n, f) = O(n^{-\beta}) \Leftrightarrow \omega_p^4(t, A^L f) = O(t^{\beta-2L})$$

*if we choose the integer  $L$  so that*

$$2L < \beta < 2L + 4.$$

**EXAMPLE.** Using our theory it is not difficult to see that if  $f(x) = (1-x^2)^\beta$  then  $E_p(n, f) = O(n^{-2\beta - (2\nu+1)/p})$ , provided that  $2\beta + (2\nu+1)/p > 0$ . We ask the reader to check this.

#### REFERENCES

- [1] BERGH, J. and LÖFSTRÖM, J., *Interpolation spaces, An Introduction*, Springer, Berlin, 1976. *MR* 58#2349.
- [2] BOCHNER, S., Sturm—Liouville and heat equation whose eigenfunctions are ultraspherical polynomials or associated Bessel functions, *Proc. Conf. Diff. Equations*, Univ. of Maryland, 1956. *MR* 18—484.
- [3] BUTZER, P. L. and BEHRENS, H., *Semi-groups of operators and approximation*, Springer-Verlag, New York, 1967. *MR* 37#5588.
- [4] DE VORE, R. and SCOTT, R., Error estimates for Gaussian quadrature and weighted  $L^1$ -polynomial approximation, Technical Report, Univ. of Wisconsin, Aug 1981.
- [5] KY, N. X., On Jackson and Bernstein type approximation theorems in the case of approximation by algebraic polynomials in  $L_p$ -spaces, *Studia Sci. Math. Hungar.* 9 (1974), 405—415. *MR* 53#6182.
- [6] LÖFSTRÖM, J. and PEETRE, J., Approximation theorems connected with generalized translations, *Math. Ann.* 181 (1969), 255—268. *MR* 40#618.
- [7] PEETRE, J. and VRETARE, L., Multiplier theorems connected with generalized translations, Technical report, Lund, 1971.
- [8] STEIN, E. M., Interpolation in polynomial classes and of Markoff's inequality, *Duke Math. J.* 24 (1957), 467—476. *MR* 19—956.
- [9] SZEGÖ, G., *Orthogonal Polynomials*, A.M.S. Coll. Publ., vol. XXIII, Providence, 4th ed., 1975. *MR* 51#8724.

(Received December 7, 1983)

## ÜBER DIE AUTOPARALLELE ABWEICHUNG VON FINSLER—OTSUKISCHEN RÄUMEN UND ANWENDUNGEN IN RÄUMEN MIT SPEZIELLEN $P$ -TENSOREN

ARTHUR MOÖR und DJERDJI F. NADJ

### § 1. Einleitung

In dem Aufsatz [5]<sup>1</sup> ist eine Übertragungstheorie definiert, die die Cartansche Theorie der Finslerräume und die der Otsukischen Räume (vgl. [7]) in sich vereinigt. In diesem verallgemeinerten Raum wollen wir in unserer vorliegenden Arbeit die Gleichung der autoparallelen Abweichung bestimmen, die bei vielen Untersuchungen — wie z. B. bei der Bestimmung des Durchmessers des Raumes, bzw. bei der Existenz der Hüllkurve der autoparallelen Linien — fundamentale Bedeutung hat (vgl. [3], [4], [8]).

In dem vom ersten Verfasser stammenden Teil I dieser Arbeit wollen wir die Gleichung der Abweichung benachbarter autoparalleler Linien in ihrer allgemeinsten Form bestimmen und für den Vektor  $P_0^i$  keine einschränkende Bedingungen voraussetzen; In dem von der Verfasserin Dj. F. Nadj geschriebenen Teil II untersuchen wir dann einige wichtige Type der  $F-O_n$ -Räume, in denen die in der Gleichung der autoparallelen Abweichung vorkommenden Grundtensoren eine besonders einfache Form haben (vgl. die Formeln (7.8a) und (7.8b)). Diese Beispiele sind durch die spezielle Form des Fundamentalvektors  $P_0^i$  der Finsler—Otsukischen  $F-O_n$ -Räume charakterisiert; wir verweisen nur auf die beiden Arbeiten [5] und [6], in denen die Übertragungsparameter  $'\Gamma_{jk}^i$  und  $''\Gamma_{jk}^i$  der durch die Relationen:  $P_0^i = l^i$  bzw.  $P_0^i = h^i$  charakterisierten Type der  $F-O_n$ -Räume bestimmt wurden, wo  $l^i$ ,  $h^i$  die von L. Berwald in [1] eingeführten Einheitsvektoren des zweidimensionalen Finslerraumes sind (vgl. [1], (2.13) und (4.5)). Wir bemerken aber schon hier, daß  $l^i$  unmittelbar auch im  $n$ -dimensionalen Fall definiert ist, sogar behält dieser Vektor auch im  $n$ -dimensionalen Raum die wichtige Eigenschaft, daß seine Cartansche kovariante Ableitung verschwindet (vgl. [2] § 7), was für  $h^i$  nur im zweidimensionalen Fall gültig ist (vgl. unsere Schlußbemerkungen in § 8).

<sup>1</sup> Vgl. die Literatur am Ende unserer Arbeit.

1980 *Mathematics Subject Classification*. Primary 53B40.

*Key words and phrases*. Local differential geometry, connections, Finsler spaces and generalizations.

## Teil I

## § 2. Grundrelationen der Finsler—Otsukischen Räume

Ein  $n$ -dimensionaler Finsler—Otsukischer  $F$ — $O_n$ -Raum ist eine Mannigfaltigkeit der Linienelemente  $(x^i, \dot{x}^i)$  in der eine Metrik durch eine in  $\dot{x}^i$  von erster Dimension homogene metrische Grundfunktion  $F(x, \dot{x})$  festgelegt ist, bzw. der metrische Grundtensor die Form

$$g_{ik}(x, \dot{x}) := \frac{1}{2} \dot{\partial}_i \dot{\partial}_k F^2(x, \dot{x}), \quad \dot{\partial}_i := \frac{\partial}{\partial \dot{x}^i}$$

hat. Dabei ist also  $F(x, \dot{x})$  die erste Grundgröße des  $F$ — $O_n$ -Raumes, für die die gewöhnlichen Regularitätsbedingungen (vgl. [8], I) bestehen sollen.

Die Übertragung der Tensoren ist durch zwei verschiedene Übertragungsparameter  $'\Gamma_{jk}^i$  und  $''\Gamma_{jk}^i$  bestimmt, wo die  $'\Gamma_{jk}^i$  bei den kontra- bzw. die  $''\Gamma_{jk}^i$  bei den kovarianten Indizes verwendet werden. Diese Theorie ist ausführlich in der Arbeit [5] in § 2 begründet, hier wollen wir nur einige wichtige und im späteren benützende Formeln angeben.

Das invariante Differential eines Tensorfeldes  $T_b^a(x, \dot{x})$  ist durch die Formeln

$$(2.1) \quad \bar{D}T_b^a = dT_b^a + (A_s^a T_b^s - A_b^s T_s^a - A_c^s T_b^a) \bar{\omega}^k(d) + \\ + ('T_s^a T_b^s - ''T_b^s T_s^a - ''T_c^s T_b^a) dx^k$$

festgelegt, wo

$$(2.2) \quad A_s^i := \frac{1}{4} F g^{ir} \dot{\partial}_r \dot{\partial}_s \dot{\partial}_k F^2 \equiv \frac{1}{2} g^{ir} g_{rs||k}, \quad (||k := F \dot{\partial}_k)$$

den Torsionstensor des Raumes bedeutet, ferner

$$\bar{\omega}^k(d) := \bar{D}l^k \equiv dl^k + '\Gamma_0^k{}_i dx^i$$

ist und die Übertragungsparameter in der Arbeit [5] § 2 durch eine längere Rechnung bestimmt sind. Schreiben wir (2.1) in der Form (vgl. [5], (2.11)—(2.13)):

$$(2.3) \quad \bar{D}T_b^a = \overset{*}{\nabla}_k T_b^a \bar{\omega}^k(d) + \overset{*}{\nabla}_k T_b^a dx^k,$$

so sind  $\overset{*}{\nabla}_k$  und  $\overset{\circ}{\nabla}_k$  eben die fundamentalen kovarianten Ableitungen des Raumes. In expliziter Form benötigen wir im folgenden nur die kovariante Ableitung

$$(2.4) \quad \overset{*}{\nabla}_k T_b^a := \partial_k T_b^a - T_b^a{}_{||s} \Gamma_0^s + '\Gamma_{sk}^a T_b^s - ''\Gamma_{bk}^s T_s^a - ''\Gamma_{ck}^s T_b^a.$$

Wir bemerken hier, daß in der kovarianten Ableitung  $\overset{*}{\nabla}_k$  bei dem Glied  $T_b^a{}_{||s}$  immer die Größe  $\Gamma_0^s$  vorkommt, während in der, im späteren benützenden kovarianten Ableitung  $\overset{*}{\nabla}_k$ , überall nur  $''\Gamma_{jk}^i$  vorhanden sein kann. Die kovariante Ableitung  $\overset{*}{\nabla}_k$  stimmt mit der in der Cartanschen Theorie vorkommenden zweiten kovarianten Ableitung überein. (Vgl. [8] IV. (1.20), bzw. [5], (2.12)).

Die Übertragungsparameter  $'\Gamma_{jk}^i$  und  $''\Gamma_{jk}^i$  sind durch die in den  $F$ — $O_n$ -Räumen immer vorausgesetzte Relation:

$$(2.5) \quad \partial_k P_j^i - P_{j||s}^i \Gamma_0^s - '\Gamma_{jk}^s P_s^i + ''\Gamma_{sk}^i P_j^s = 0$$



miteinander verbunden, wo  $P_j^i(x, \dot{x})$  neben der Grundfunktion  $F(x, \dot{x})$  die zweite Grundgröße des  $F-O_n$ -Raumes ist. Es soll immer  $\text{Det}(P_j^i) \neq 0$  sein, wodurch die Existenz des inversen Tensors:  $Q_j^i$  gesichert ist.

Für die Bestimmung der beiden Übertragungsparameter hat man noch außer (2.5) die Relation

$$(2.6) \quad \nabla_k g_{ij} \equiv \partial_k g_{ij} - 2A_{ijs}' \Gamma_{0k}^s - \Gamma_{ijk} - \Gamma_{jik} = 0,$$

woraus mit (2.5) zusammen, die in  $(j, k)$  symmetrischen  $\Gamma_{jk}^i$  und die  $\Gamma_{jk}^i$  bestimmt werden können (vgl. [5], § 2; insbesondere (2.14) und die nachfolgenden Zeilen).

Wir bemerken noch, daß alle Größen der  $F-O_n$ -Räume — außer der Grundfunktion  $F(x, \dot{x})$  — in den  $\dot{x}^i$  homogen von nullter Dimension sind und nach (2.2) auch  $A_{0ik} = A_{i0k} = A_{ik0} = 0$  ist.

### § 3. Gleichung der autoparallelen Abweichung

Die autoparallelen Kurven sind die Lösungskurven des Differentialgleichungssystems

$$(3.1) \quad \frac{\bar{D}}{ds} \frac{dx^i}{ds} \equiv \frac{d^2 x^i}{ds^2} + \Gamma_{jk}^i \left( x, \frac{dx}{ds} \right) \frac{dx^j}{ds} \frac{dx^k}{ds} = 0,$$

wo die Kurve die Gesamtheit  $(x^i(s), l^i(s))$  der tangentialen Linienelemente ist, d. h. es gilt jetzt und im folgenden immer  $l^i(s) = dx^i/ds$ . Der Parameter „ $s$ “ in (3.1) bedeutet entweder die Bogenlänge, oder aber einen affinen Parameter bezüglich der affinen Übertragung  $\Gamma_{jk}^i$  (vgl. [9], § 2). Die zur autoparallelen Kurve (3.1) unendlich benachbarte autoparallele Kurve, die von einem Punkt  $O$  von (3.1) ausgeht, habe die charakteristische Gleichung:

$$(3.2) \quad \frac{\bar{D}}{d\sigma} \frac{d\psi^i}{d\sigma} \equiv \frac{d^2 \psi^i}{d\sigma^2} + \Gamma_{jk}^i \left( \psi, \frac{d\psi}{d\sigma} \right) \frac{d\psi^j}{d\sigma} \frac{d\psi^k}{d\sigma} = 0,$$

wo  $\sigma$  die Bogenlänge, oder möglicherweise den affinen Parameter der durch (3.2) bestimmten Kurve bedeutet. Da die beiden durch (3.1) und (3.2) bestimmten Kurven unendlich benachbarte Kurven sein sollen, gilt:

$$(3.3a) \quad \psi^i(\sigma) = x^i(s) + \xi^i(s),$$

$$(3.3b) \quad \frac{d\sigma}{ds} = 1 + \varepsilon(s),$$

wo  $\xi^i(s)$  einen infinitesimalen Vektor bzw.  $\varepsilon(s)$  einen infinitesimalen Skalar bedeutet. Für eine eingehendere Begründung vgl. [3] und [7].

Mit derselben Methode, die in [3] verwendet wurde und die auf Grund von (3.3a) und (3.3b) in der Vernachlässigung der Glieder höherer Größenordnung in  $\xi^i(s)$  und  $\varepsilon(s)$  in der Gleichung (3.2) besteht, ferner wenn  $\frac{d\xi^i}{ds}$  bzw.  $\frac{d}{ds} \frac{D\xi^i}{ds}$  immer durch

$$\frac{d\xi^i}{ds} = \frac{\bar{D}\xi^i}{ds} - \Gamma_{j0}^i \xi^j \quad \text{bzw.} \quad \frac{d}{ds} \frac{\bar{D}\xi^i}{ds} = \frac{\bar{D}^2 \xi^i}{ds^2} - \Gamma_{j0}^i \frac{\bar{D}\xi^j}{ds}$$

ersetzt wird, so erhält man als die der Gleichung (3.7) von [3] entsprechende Gleichung der autoparallelen Abweichung im  $F-O_n$ -Raum:

$$(3.4) \quad \frac{|\bar{D}^2 \xi^i}{ds^2} - \frac{dx^i}{ds} \frac{d\varepsilon}{ds} + \hat{F}^i_j(x, \dot{x}) \frac{\bar{D} \xi^j}{ds} + \hat{R}^i_j(x, \dot{x}) \xi^j = 0,$$

wo, unter Verwendung der Schoutenschen Symbolik:

$$(3.5) \quad \hat{F}^i_j(x, \dot{x}) := 2' \Gamma_{[0]j}^i + ' \Gamma_{r[k]j}^i l^r l^k$$

$$(3.6) \quad \hat{R}^i_j(x, \dot{x}) := 'R_{0j}^i + 2' \Gamma_{r[k]j}^i l^r l^k + 2(' \nabla_0 ' \Gamma_{[r]j}^i) l^r$$

und  $'R_{abj}^i$  den durch  $' \Gamma_{jk}^i$  bestimmten Hauptkrümmungstensor, d. h.

$$(3.7) \quad 'R_{abj}^i := \partial_j ' \Gamma_{ab}^i - \partial_b ' \Gamma_{aj}^i - ' \Gamma_{ab||r}^i ' \Gamma_0^r j + ' \Gamma_{a||r}^i ' \Gamma_0^r b + ' \Gamma_{ab}^i ' \Gamma_{rj}^i - ' \Gamma_{aj}^i ' \Gamma_{rb}^i$$

bedeutet.

Die bisherigen Resultate fassen wir im folgenden Satz zusammen.

**SATZ 1.** *Die allgemeinste Form der autoparallelen Abweichung in den Finsler-Otsukischen Räumen bezüglich der durch (3.1) charakterisierten autoparallelen Kurven ist durch (3.4) angegeben, wo die Koeffizienten von  $\frac{\bar{D} \xi^i}{ds}$  und  $\xi^i$  durch (3.5) und (3.6) angegeben sind.*

#### § 4. Der Typ $' \Gamma_{0k}^i = '' \Gamma_{0k}^i + l^i q_k$

In dem Aufsatz [5] im Satz 1 haben wir eine Übertragung der  $F-O_n$ -Räume bestimmt, die mit der Cartanschen Übertragung der Finsler-Räume in einem sehr engen Zusammenhang war. Jetzt wollen wir in diesem Paragraphen den im Aufsatz [5] im Satz 1 entwickelten Typ verallgemeinern und von einer anderen Seite ausgehend charakterisieren. Es besteht der

**SATZ 2.** *Besteht die Relation*

$$(4.1) \quad ' \Gamma_{0k}^i = '' \Gamma_{0k}^i + l^i q_k,$$

wo  $q_k$  einen kovarianten Vektor bedeutet, so ist  $'' \Gamma_{jk}^i$  mit dem Cartanschen Übertragungsparameter  $\Gamma_{jk}^{*i}$  identisch und es gelten die Relationen

$$(4.2a) \quad ' \Gamma_{jk}^i = '' \Gamma_{jk}^i + Q_m^i '' \nabla_k P_j^m,$$

$$(4.2b) \quad '' \nabla_k P_0^i = P_0^i q_k,$$

wo  $'' \nabla_k$  die allein mit  $'' \Gamma_{jk}^i$  gebildete kovariante Ableitung — im Fall  $'' \Gamma = \Gamma^*$  die Cartansche kovariante Ableitung — bezeichnet. (Den Fall  $q_k \equiv 0$  siehe in [5], (2.21).)

**BEWEIS.** Die erste Hälfte des Satzes ist fast trivial, da wenn  $' \Gamma_{0k}^i$  aus (4.1) in die Relation (2.6) substituiert wird, so entsteht für  $'' \Gamma_{jk}^i$  eben jenes Gleichungssystem, die zur Bestimmung der in  $(j, k)$  symmetrischen Übertragungsparameter  $\Gamma_{jk}^{*i}$  dient (vgl. [8], Kap. III. (1.23)).

Betreffs der Relation (4.2a) des Satzes setzen wir in die in den  $F-O_n$ -Räumen fundamentale Gleichung (2.5) statt  $' \Gamma_{0k}^i$  die entsprechende Form von (4.1) ein,

womit in Hinsicht auf die Homogenität von nullter Dimension von  $P_j^i(x, \dot{x})$  in den  $\dot{x}^i$  die Relation

$$(4.3) \quad \partial_k P_j^i - P_{j||s}^i {}''\Gamma_{0k}^s - {}'\Gamma_{jk}^s P_s^i + {}''\Gamma_{sk}^i P_j^s = 0$$

entsteht. Addieren und subtrahieren wir zur linken Seite  ${}''\Gamma_{jk}^s P_s^i$ , dann erhalten wir nach einer Kontraktion mit dem Tensor  $Q_i^m$ , nach einigen Indexvertauschungen und auf Grund von  $Q_i^m P_s^i = \delta_s^m$  eben (4.2a).

Endlich überschieben wir noch (4.3) mit  $l^j$ , so wird nach (4.1)

$$(4.3^*) \quad \begin{aligned} &\partial_k P_0^i - P_s^i \partial_k l^s - P_{0||s}^i {}''\Gamma_{0k}^s + P_j^i l_{||s}^j {}''\Gamma_{0k}^s - \\ &- ({}''\Gamma_{0k}^s + l^s q_k) P_s^i + {}''\Gamma_{sk}^i P_0^s = 0 \end{aligned}$$

entstehen. Wir berechnen nun  $\partial_k l^s$  und  $l_{||s}^j$ . Unter Beachtung von (2.6) und (2.2) werden:

$$(4.4) \quad \partial_k l^s = \partial_k (g_{ab} \dot{x}^a \dot{x}^b)^{-1/2} \dot{x}^s \equiv -{}''\Gamma_{00k}^s l^s, \quad l_{||s}^j = \delta_s^j - l^j l_s$$

gelten, wodurch (4.3\*) in

$$(4.5) \quad \partial_k P_0^i - P_{0||s}^i {}''\Gamma_{0k}^s + {}''\Gamma_{sk}^i P_0^s = P_0^i q_k$$

übergeht und das stimmt mit (4.2b) überein.

Dieser Satz, die für die Bestimmung der Übertragungsparameter bzw. der Tensoren (3.5) und (3.6) eine sehr wichtige Rolle spielt, kann nicht vollständig umgekehrt werden. Es gilt aber der

**SATZ 3.** Besteht die Relation (4.2b), wo  $q_k$  einen kovarianten Vektor bezeichnet, ist ferner

$$(4.6) \quad P_0^i Q^*_{i^j} = l^j$$

wo  $Q^*_{i^j}$  den inversen Tensor von  $(P_s^i + l^b P_{b||s}^i)$  bedeutet (vgl. [5], (2.17)), so besteht auch (4.1).

**BEWEIS.** Überschieben wir (2.5) mit  $l^j$ , beachten wir dann die beiden Identitäten von (4.4) und subtrahieren wir die erhaltene Relation aus (4.2b) bzw. aus der mit (4.2b) übereinstimmenden Gleichung (4.5), so wird

$$(P_0^i l_s + P_{0||s}^i)({}'\Gamma_{0k}^s - {}''\Gamma_{0k}^s) = P_0^i q_k.$$

Beachten wir jetzt, daß  $P_{0||s}^i = l^b P_{b||s}^i + P_b^i l_{||s}^b$  ist, so wird nach der zweiten Formeln von (4.4):

$$(P_s^i + l^b P_{b||s}^i)({}'\Gamma_{0k}^s - {}''\Gamma_{0k}^s) = P_0^i q_k,$$

woraus nach einer Kontraktion mit  $Q^*_{i^j}$  und in Hinsicht auf die Bedingung (4.6) unmittelbar (4.1) entsteht, w. z. b. w.

**BEMERKUNG.** Offenbar bleibt der Satz im wesentlichen gültig, falls statt (4.6) die Relation:

$$P_0^i Q^*_{i^j} = \lambda l^j, \quad \lambda = \lambda(x, \dot{x}), \quad (\lambda = \text{Skalar})$$

besteht, nur  $q_k$  verändert sich mit dem skalaren Faktor  $\lambda$ .

## Teil II

§ 5. Der Typ  $P_0^i = \lambda l^i$  der  $F-O_n$ -Räume

Im Aufsatz [5] haben wir den Fall  $P_0^i = l^i$  ausführlich untersucht, jetzt wollen wir den Typ  $P_0^i = \lambda(x, \dot{x}) l^i$  eingehend untersuchen. Es wird sich zeigen, daß jetzt eben die Relation (4.1) mit  $q_k = \lambda_{|k} \lambda^{-1}$  entsteht, wo „ $|k$ “ die Cartansche kovariante Ableitung bezeichnet.

BEMERKUNG. Aus dem Satz 2 folgt, daß die Cartansche kovariante Ableitung, falls (4.1) besteht, mit  ${}''\nabla_k$  übereinstimmt.

Es zeigt sich in den folgenden Untersuchungen, daß die Type der  $F-O_n$ -Räume, in denen (4.1) gültig ist, z. B. durch die ziemlich einfache Bedingung  $P_0^i = \lambda l^i$  realisierbar sind. Es gilt der

SATZ 4. Ist in einem  $F-O_n$ -Raum  $P_0^i = \lambda l^i$ , wo  $\lambda = \lambda(x, \dot{x})$  einen Skalar bedeutet, der in den  $\dot{x}^i$  homogen von nullter Dimension ist, so ist (4.1) gültig, genauer, es ist

$$(5.1) \quad {}'\Gamma_0^j{}^k = {}''\Gamma_0^j{}^k + l^j {}''\nabla_k \ln \lambda,$$

wo in diesem Fall  ${}''\nabla_k$  die Cartansche kovariante Ableitung bestimmt;  ${}''\Gamma_j^i{}^k$  sind die Cartanschen Übertragungsparameter, ferner es gilt<sup>2</sup>

$$(5.2) \quad {}''\nabla_k P_0^i = l^i {}''\nabla_k \lambda.$$

BEWEIS. Die Formel

$$(5.3) \quad {}'\Gamma_0^s{}^k (P_s^i + l^b P_{b||s}^i) = l^b \partial_k P_b^i + {}'\Gamma_s^i{}^k P_0^s$$

ist in allen  $F-O_n$ -Räumen gültig, da die Formel (5.3) eine unmittelbare Folgerung von (2.5) ist (vgl. [5], (2.16)). Ist nun  $P_0^i = \lambda l^i$ , so ist

$$l^b P_{b||s}^i = P_{0||s}^i - P_b^i l_{||s}^b = P_{0||s}^i - P_b^i (\delta_s^b - l^b l_s),$$

d. h. es wird in diesem Fall

$$(5.4) \quad P_s^i + l^b P_{b||s}^i = \lambda \delta_s^i + \lambda_{||s} l^i.$$

Durch eine unmittelbare Rechnung kann gezeigt werden, daß der inverse Tensor von (5.4), d. h.  $Q_i^{*j}$  in diesem Fall die Form:

$$(5.5) \quad Q_i^{*j} = \lambda^{-1} \delta_i^j - \lambda^{-2} \lambda_{||i} l^j$$

hat.

Wir umformen noch die rechte Seite von (5.3). Es wird:

$$l^b \partial_k P_b^i = \frac{1}{F} \dot{x}^b \partial_k P_b^i = \frac{1}{F} \partial_k P_b^i \dot{x}^b = \frac{1}{F} \partial_k (F P_0^i) = \frac{1}{F} \partial_k \lambda \dot{x}^i,$$

d. h. es ist

$$(5.6) \quad l^b \partial_k P_b^i = (\partial_k \lambda) l^i.$$

<sup>2</sup> Unser Satz 4 ist also die unmittelbare Verallgemeinerung des Satzes 1 von [5].

Setzt man jetzt die entsprechenden Größen von (5.4) und (5.6) in (5.3) ein, so wird nach einer Kontraktion mit  $Q_i^{*j}$  wegen  $\lambda_{||i} l^i \equiv \lambda_{||0} = 0$

$${}^i\Gamma_{0k}^j = (\partial_k \ln \lambda) l^j + {}^i\Gamma_{0k}^j - l^j (\ln \lambda)_{||i} {}^i\Gamma_{0k}^i$$

und das ist eben die Formel (5.1).

Da (5.1) dieselbe Form wie (4.1) hat, wird nach dem Satz 2  ${}^i\Gamma_{jk}^i$  mit den Cartanschen Übertragungsparametern übereinstimmen<sup>3</sup>, ferner es folgt noch nach (4.2b) wegen  ${}^i\nabla_k l^i \equiv l_{||k}^i = 0$  die Formel (5.2), womit der Satz 4 vollständig bewiesen ist.

## § 6. Der Typ $P_0^i = \mu h^i$ der $F-O_2$ -Räume

In diesem Paragraphen wollen wir den Typ  $P_0^i = \mu h^i$  der 2-dimensionalen  $F-O_n$ -Räume untersuchen, wo  $\mu(x, \dot{x})$  einen Skalar bedeutet, der in  $\dot{x}^i$  homogen von nullter Dimension ist;  $h^i$  bedeutet der Berwaldsche — auf  $l^i$  orthogonale — Einheitsvektor:

$$(6.1) \quad h^i = -\varepsilon^{ij} l_j, \quad \varepsilon^{ij} := \begin{pmatrix} 0 & g^{1/2} \\ -g^{-1/2} & 0 \end{pmatrix} \quad g := \text{Det}(g_{ik})$$

(vgl. [1], (4.1)–(4.5)). Dieser Typ realisiert auch (4.1) — wie das gezeigt wird — zwar nur im 2-dimensionalen Fall. Es besteht der

**SATZ 5.** *Ist in einem  $F-O_n$ -Raum  $P_0^i = \mu h^i$ , wo  $\mu(x, \dot{x})$  einen in  $\dot{x}^i$  von nullter Dimension homogenen Skalar bezeichnet, so besteht die Relation (4.1), genauer, es ist*

$$(6.2) \quad {}^i\Gamma_{0k}^i = {}^i\Gamma_{0k}^i + l^i {}^i\nabla_k \ln \mu.$$

Es besteht auch die Formel

$$(6.2a) \quad {}^i\nabla_k P_0^i = h^i \nabla_k \mu$$

und die Übertragungsparameter  ${}^i\Gamma_{jk}^i$  sind mit den Cartanschen  $\Gamma_{jk}^{*i}$  Übertragungsparametern identisch.

**BEWEIS.** Wir beginnen den Beweis mit der Berechnung der entsprechenden Form der Gleichung (5.3). Die Gleichung (5.3) ist selbstverständlich in allen  $F-O_n$ -Räumen gültig. Nach  $P_0^i = \mu h^i$  wird:

$$l^b P_{b||s}^i = P_{0||s}^i - (\delta_s^b - l^b l_s) P_b^i$$

d. h.

$$P_s^i + l^b P_{b||s}^i = \mu_{||s} h^i + \mu h_{||s}^i + \mu h^i l_s.$$

Im 2-dimensionalen Finslerraum ist nun (vgl. [1], (6.2)–(6.4), bzw. [6], (1.13)):

$$\mu_{||s} = \mu_{\mathfrak{g}} h_s, \quad h_{||s}^i = -(l^i + \sqrt{A} h^i) h_s,$$

wo der Index  $\mathfrak{g}$  die Operation:  $\mu_{\mathfrak{g}} := \mu_{||k} h^k$  bedeutet und nach (2.2)  $A := A' A$ ,

<sup>3</sup> Vgl. [8], Kap. III. (1.23).

mit  $A^t := A^{t,k}$  ist. Somit wird

$$(6.3) \quad P_s^i + l^b P_{b||s}^i = \mu \left\{ \left[ \left( \frac{\mu_s}{\mu} - \sqrt{A} \right) h_s + l_s \right] h^i - l^i h_s^i \right\}.$$

Auf Grund einer unmittelbaren Rechnung kann gezeigt werden, daß der inverse Tensor von (6.3)

$$(6.4) \quad Q^{*t}_i = \mu^{-1} \left\{ l^t h_i + \left[ \left( \frac{\mu_s}{\mu} - \sqrt{A} \right) l^t - h^t \right] l_i \right\}$$

ist, nur müssen wir die Identität  $l^t l_s + h^t h_s \equiv \delta_s^t$  beachten.

Wir berechnen jetzt die rechte Seite von (5.3). Es ist nach (4.4)

$$l^b \partial_k P_b^i = \partial_k (\mu h^i) - P_b^i \partial_k l^b = (\partial_k \mu) h^i + \mu \partial_k h^i + \mu h^i{}'' \Gamma_{00k}.$$

Nach der Formel (1.15) von [6] ist

$$\partial_k h^i = -(\sqrt{A} h_t{}' \Gamma_{0k}^t h^i + {}''\Gamma_{tk}^i h^t + {}''\Gamma_{0k}^t h_t l^i),$$

was man übrigens aus (6.1) berechnen kann, wenn man die entsprechenden Formeln von [1] § 4, und für  $\partial_k g_{ij}$  unsere Formel (2.6) verwendet. Nach unseren letzten beiden Formeln wird somit die rechte Seite von (5.3)

$$(6.5) \quad l^b \partial_k P_b^i + {}''\Gamma_{sk}^i P_0^s = \mu (\partial_k \ln \mu - \sqrt{A} h_s{}' \Gamma_{0k}^s + {}''\Gamma_{00k}) h^i - \mu h_s {}''\Gamma_{0k}^s l^i.$$

Substituieren wir jetzt (6.3) und (6.5) in die Formel (5.3), so wird nach einer Kontraktion mit dem in (6.4) angegebenen Tensor  $Q^{*t}_i$ :

$${}'\Gamma_{0k}^t = l^t \{ \partial_k \ln \mu - (\ln \mu)_{||r} h^r h_s {}''\Gamma_{0k}^s - \sqrt{A} h_s ({}'\Gamma_{0k}^s - {}''\Gamma_{0k}^s) \} + {}''\Gamma_{0k}^s (l^t l_s + h^t h_s).$$

Beachten wir jetzt zweimal die schon verwandte Identität  $h^t l_s = \delta_s^t - l^t l_s$  und die Homogenität von nullter Dimension in  $\dot{x}^i$  von  $\mu$ , d. h.  $\mu_{||0} = 0$ , so wird

$${}'\Gamma_{0k}^t = {}''\Gamma_{0k}^t - \sqrt{A} h_s ({}'\Gamma_{0k}^s - {}''\Gamma_{0k}^s) l^t + l^t {}''\nabla_k \ln \mu,$$

was offensichtlich in der Form

$$({}'\Gamma_{0k}^s - {}''\Gamma_{0k}^s) (\delta_s^t + \sqrt{A} h_s l^t) = l^t {}''\nabla_k \ln \mu$$

geschrieben werden kann. Nach einer Kontraktion von beiden Seiten mit  $(\delta_i^t - \sqrt{A} h_t l^i)$  bekommt man die Formel (6.2).

Da (6.2) die Form von (4.1) hat und jetzt  $q_k = {}''\nabla_k \ln \mu$  ist, kann der Satz 2 verwendet werden. Beachten wir die Formel  $P_0^t = \mu h^t$ , so gibt (4.2b) wegen der Form von  $q_k$  unmittelbar (6.2a). Aus dem Satz 2 folgt noch, daß auch jetzt  ${}''\Gamma_{jk}^i$  bzw.  ${}''\nabla_k$  die Cartansche Übertragungsparameter bzw. die Cartansche kovariante Ableitung bestimmen.

Der Satz 5 drückt aus, daß der durch  $P_0^i = \mu(x, \dot{x}) h^i$  gekennzeichnete Typ der  $F - O_n$ -Räume zu dem Typ  $P_0^i = \lambda(x, \dot{x}) l^i$  ähnlich ist, nur  ${}''\nabla_k P_0^i$  ist nach (5.2) und (6.2a) in den beiden Typen voneinander verschieden.



## § 7. Sätze über die Gleichung der autoparallelen Abweichung

Die im vorigen betrachteten Type waren alle dadurch charakterisiert, daß der Zusammenhang von  $\Gamma_{0k}^i$  und  $\Gamma_{jk}^i$  durch eine Relation von der Form (4.1) bestimmt war (vgl. (4.1), (5.1) und (6.2)). Nach Satz 2 ist in diesen Räumen  $\Gamma_{jk}^i = \Gamma_{jk}^{*i}$ , wo die  $\Gamma_{jk}^{*i}$  die Cartanschen Übertragungsparameter bedeuten.

Setzen wir nun

$$(7.1) \quad \Gamma_{jk}^i = \Gamma_{jk}^{*i} + A_j^i,$$

so werden im allgemeinen die Tensoren (3.5) und (3.6) auch ein tensorielles von  $A_j^i$  abhängiges Glied enthalten, welches die Abweichung von der Cartanschen Theorie bestimmt. Wegen der Symmetrie von  $\Gamma_{jk}^{*i}$  in  $(j, k)$  ist offensichtlich  $\Gamma_{[j]k}^i = A_{[j]k}^i$ , und nach (3.7)

$$(7.2) \quad R_{abj}^i = K_{abj}^i + \Phi_{abj}^i,$$

wo  $K_{abj}^i$  den aus  $\Gamma_{ab}^{*i}$  gebildeten *Hauptkrümmungstensor* des Raumes (vgl. [8], Kap. IV. (1.7)) bzw.  $\Phi_{abj}^i$  den Tensor

$$(7.3) \quad \Phi_{abj}^i := A_{ab|j}^i - (\Gamma_{ab||j}^{*i} + A_{ab||j}^i) A_0^i + A_{ab}^i A_{||j}^i - [j|b]$$

bedeutet. Dabei ist „ $|j$ “ wieder die mit  $\Gamma_{ab}^{*i}$  gebildete, d.h. die *Cartansche kovariante Ableitung* und  $[j|b]$  bedeutet den ganzen vorigen Ausdruck mit vertauschten Indizes  $j, b$ . Offenbar kommen in den im vorigen untersuchten Typen, d.h. in §§ 4–6 direkt diese Tensoren vor.

Nach (3.5) und (7.1) ist nun

$$(7.4) \quad \hat{F}_j^i = A_{ab||j}^i l^a l^b + A_0^i j - A_j^i,$$

da (vgl. z. B. [2], (44)):

$$(7.5) \quad \Gamma_{ab||j}^{*i} l^a = A_j^i l_b, \quad \Gamma_{ab||j}^{*i} l^a l^b = 0.$$

Unter Beachtung der zweiten Relation von (4.4) wird aus (7.4)

$$(7.6) \quad \hat{F}_j^i = A_{00||j}^i - 2A_{j0}^i + 2A_0^i l_j.$$

Wir bestimmen jetzt  $\hat{R}_j^i$  aus (3.6) in Hinsicht auf (7.2) und (7.3). Unter Beachtung der Relationen (7.5) wird:

$$\begin{aligned} \hat{R}_j^i &= K_{00j}^i + A_{00||j}^i - A_{0j|0}^i + A_j^i A_{00}^i + A_{r||j}^i l^r A_{00}^i - \\ &\quad - A_{r||j}^i l^r l^k A_{j0}^i + A_{00}^i A_{||j}^i - A_{0j}^i A_{||0}^i + 2\nabla_0 A_{[0j]}^i - 2A_{[rj]}^i \nabla_0 l^r. \end{aligned}$$

Nun ist wegen  $F_{||j}^i = 0$  nach (7.1)

$$\nabla_0 l^r = l^r A_{000},$$

somit wird

$$(7.7) \quad \begin{aligned} \hat{R}_j^i &= K_{00j}^i + A_{00||j}^i - A_{0j|0}^i + A_j^i A_{||0}^i + A_{0j||}^i A_{00}^i - \\ &\quad - A_{00||j}^i A_j^i - A_{0j}^i A_{||0}^i + 2A_{(0j)}^i A_j^i - 2A_{00}^i A_{j00} + 2\nabla_0 A_{[0j]}^i + A_j^i A_{000}. \end{aligned}$$

Für die autoparallele Abweichung (3.4) gelten auf Grund von (7.6) und (7.7) die folgenden Sätze:



SATZ 6. In einem durch  $P_0^i = \lambda l^i$ ,  $\lambda_{|k} = 0$  charakterisierten  $F-O_n$ -Raum bestehen für die Tensoren (7.6) und (7.7) der autoparallelen Abweichung:

$$(7.8a) \quad \hat{F}_j^i = -2Q_s^i P_{j|0}^s,$$

$$(7.8b) \quad \hat{R}_j^i = K_{00j}^i - Q_{r|0}^i P_{j|0}^r - \nabla_0(Q_s^i P_{j|0}^s),$$

wenn das Symbol „ $_{|k}$ “ die Cartansche kovariante Ableitung bedeutet. Ist  $P_{j|0}^m = 0$ , so stimmt die Gleichung (3.4) der autoparallelen Abweichung mit der der Finslerräume überein.

SATZ 7. In einem durch  $P_0^i = \mu h^i$ ,  $\mu_{|k} = 0$  charakterisierten  $F-O_2$ -Raum bestehen für die Tensoren (7.6) und (7.7) der autoparallelen Abweichung wieder die Relationen (7.8a) und (7.8b) und bezüglich der Relation  $P_{j|0}^m = 0$  verhalten sich diese Räume zu denen, die im vorigen Satz gekennzeichnet wurden, vollständig analog.

BEWEIS DER SÄTZE 6 UND 7. Nach den Sätzen 4 und 5 ist  ${}''\Gamma_{jk}^i = \Gamma_j^{*i}{}_k$ , d. h. die  ${}''\Gamma_{jk}^i$  bestimmen die Cartanschen Übertragungsparameter. In beiden Sätzen hat  ${}'\Gamma_{0k}^i$  die Form (4.1), wie das nach den Sätzen 4 und 5 unmittelbar verifiziert werden kann. Es kann somit auf  ${}'\Gamma_{0k}^i$  der allgemeine Satz 2 verwendet werden;  ${}'\Gamma_{jk}^i$  hat also die Form (4.2a), wo  ${}''\nabla_k$  in die Cartansche, d. h. mit  $\Gamma_j^{*i}{}_k$  gebildete kovariante Ableitung: „ $_{|k}$ “ übergehen wird. Es ist also

$${}'\Gamma_{jk}^i = \Gamma_j^{*i}{}_k + Q_m^i P_{j|k}^m$$

Vergleichen wir das mit (7.1), so sieht man, daß in (7.1)

$$(7.9) \quad A_{jk}^i = Q_m^i P_{j|k}^m$$

besteht, und wegen  $l_{|k}^j = 0$ ,  $h_{|k}^j = 0$  (vgl. [1], (2.13) (c) und (4.5)), wird:

$$(7.9a) \quad A_{0k}^i = Q_m^i l^m \lambda_{|k} \quad \text{bzw.} \quad A_{0k}^i = Q_m^i h^m \mu_{|k}.$$

Wegen unserer Annahmen über  $\lambda_{|k}$  bzw.  $\mu_{|k}$  folgt somit nach (7.6) und (7.9a), daß die Formel (7.8a) in beiden Sätzen gültig ist, da nach (7.9a)  $A_{0k}^i = 0$  bestehen muß. Auf Grund der Relation (7.9) bekommt man aus (7.7):

$$\hat{R}_j^i = K_{00j}^i + A_{i0}^i A_{j0}^i - \nabla_0 A_{j0}^i \equiv K_{00j}^i + Q_m^i P_{i|0}^m Q_r^i P_{j|0}^r - \nabla_0(Q_m^i P_{j|0}^m).$$

Wir müssen noch nur  $Q_m^i P_i^m = \delta_i^i$  beachten, dann bekommt man unmittelbar aus der letzten Gleichung (7.8b), da die Leibnizsche Regel auf die Cartansche kovariante Ableitung gültig ist.

Ist noch  $P_{j|0}^m = 0$ , so folgt nach (7.8a) und (7.8b) unmittelbar  $\hat{F}_j^i = 0$ ,  $\hat{R}_j^i = K_{00j}^i$ ; die Gleichung der autoparallelen Abweichung, d. h. die Gleichung (3.4) geht also in die entsprechende Formel der Finslerräume über, womit die Sätze 6 und 7 vollständig bewiesen sind.

Wir bemerken noch, daß wegen  $P_j^m Q_m^i = \delta_j^i$  in (7.9a)  $Q_m^i l^m = \lambda^{-1} l^i$  bzw.  $Q_m^i h^m = \mu^{-1} l^i$  ist, wie das leicht bestätigt werden kann. Aus (7.9) folgt also auch (5.1) bzw. (6.2).

# § 8. Schlußbemerkungen

Die Sätze 6 und 7 behandeln zwei Räume, die einen sehr ähnlichen Charakter haben. Doch ist zwischen diese Type der  $F-O_n$ -Räume einen wesentlichen Unterschied, denn der Vektor  $l^i$  ist auch in den  $n$ -dimensionalen Räumen definiert, während der Berwaldsche Vektor  $h^i$  nur in den zweidimensionalen Raum eingeführt wurde. Man könnte aber durch die Formel

$$h^i = (A_i A^i)^{-1/2} A^i$$

den Vektor  $h^i$  auch in den  $n$ -dimensionalen Räumen ( $n > 2$ ) definieren, doch würden dann mehrere günstige Eigenschaften von  $h^i$  verloren gehen. Z. B. im  $n$ -dimensionalen Fall ( $n > 2$ ) werden die Relationen

$$\mu_{||s} = \mu_s h_s, \quad h_{||s}^i = -(l^i + \sqrt{A} h^i) h_s, \quad h_{||k}^i = 0$$

im allgemeinen nicht bestehen, zwar der Satz 1 der autoparallelen Abweichung und auch die Sätze 2—4, die sich auf einzelne Type der  $F-O_n$ -Räume beziehen, auch im  $n$ -dimensionalen Fall gültig sind.

Der Satz 5 wird aber im  $n$ -dimensionalen Fall — vermutlich — eine wesentlich andere Form haben.

# LITERATURVERZEICHNIS

- [1] BERWALD, L., On Finsler and Cartan geometries III, *Ann. of Math.* **42** (1941), 84—112 MR 2—304.
- [2] CARTAN, E., *Les espaces de Finsler*, Actualités scientifiques et industrielles, No 79, Hermann et Cie, Paris, 1934.
- [3] MOÓR, A., Über die autoparallele Abweichung in allgemeinen metrischen Linienelementräumen, *Publ. Math. Debrecen* **5** (1957), 102—118. MR 19—980.
- [4] MOÓR, A., Über verschiedene geodätische Abweichungen in Weyl—Otsukischen Räumen, *Publ. Math. Debrecen* **28** (1981), 247—258. MR 83c: 53031.
- [5] MOÓR, A., Über die Begründung von Finsler—Otsukischen Räumen und ihre Dualität, *Tensor N. S.* **37** (1982), 121—129.
- [6] NADJ, D. F., On special two-dimensional Finsler—Otsuki space, *Rev. Res. Fac. Sci. Univ. Novi Sad, M.S.* **14** (1984), no. 2, 135—146.
- [7] OTSUKI, T., On general connections I, *Math. J. Okayama Univ.* **9** (1959—60), 99—164. MR 22# 2954.
- [8] RUND, H., Eine Krümmungstheorie der Finslerschen Räume, *Math. Ann.* **125** (1952), 1—18. MR 14—499.
- [9] RUND, H., *The differential geometry of Finsler spaces*, Springer-Verlag, Berlin—Göttingen—Heidelberg, 1959. MR 21# 4462.
- [10] VARGA, O., Über affinzusammenhängenden Mannigfaltigkeiten von Linienelementen insbesondere deren Äquivalenz, *Publ. Math. Debrecen* **1** (1949), 7—17. MR 11—134.

(Eingegangen am 19. Dezember 1983)

ERDÉSZETI ÉS FAIPARI EGYETEM  
MATEMATIKA TANSZÉK  
POSTAFIÓK 132  
H—9401 SOPRON  
HUNGARY



# SEPARABILITY IN THE UNIFORM TOPOLOGY

LIAQAT ALI KHAN

Let  $X$  be a non-empty completely regular Hausdorff space,  $E$  a Hausdorff topological vector space (TVS), and  $C_b(X, E)$  the vector space of all continuous bounded  $E$ -valued functions on  $X$ . Let  $C_0(X, E)$  denote the subspace of  $C_b(X, E)$  consisting of those functions which vanish at infinity. When  $E$  is the real field, these spaces are denoted by  $C_b(X)$  and  $C_0(X)$ . We shall denote by  $C_b(X) \otimes E$  the vector subspace spanned by the set of all functions of the form  $g \otimes a$ , where  $g \in C_b(X)$ ,  $a \in E$ , and  $(g \otimes a)(x) = g(x)a$  ( $x \in X$ ). The *uniform topology*  $\mathcal{U}$  on  $C_b(X, E)$  is the linear topology which has a base of neighbourhoods of 0 consisting of all sets of the form  $\{f \in C_b(X, E) : f(X) \subseteq W\}$ , where  $W$  is a neighbourhood of 0 in  $E$ . Note that, on  $C_b(X)$ ,  $\mathcal{U}$  is the norm topology given by  $\|f\| = \sup \{|f(x)| : x \in X\}$ .

The following fundamental results were obtained by M. and S. Krein (see [4]) and Semadeni and Zbijiński [1]. In ([3], Theorems 2.3 and 2.4), Summers has given their alternate proofs.

**THEOREM 1.**  $(C_b(X), \|\cdot\|)$  is separable iff  $X$  is a compact metric space.

**THEOREM 2.** If  $X$  is locally compact, then  $(C_0(X), \|\cdot\|)$  is separable iff  $X$  is a  $\sigma$ -compact metric space.

The purpose of this note is to extend these characterizations of separability to the case of vector-valued functions. In the sequel,  $E$  will denote a real Hausdorff TVS with non-trivial dual  $E'$ .

**THEOREM 3.**  $(C_b(X) \otimes E, \mathcal{U})$  is separable iff  $X$  is a compact metric space and  $E$  is separable.

**PROOF.** Suppose  $(C_b(X) \otimes E, \mathcal{U})$  is separable. Let  $\varphi \in E'$  with  $\varphi \neq 0$ . Define  $T_\varphi : (C_b(X) \otimes E, \mathcal{U}) \rightarrow (C_b(X), \|\cdot\|)$  by  $T_\varphi(f) = \varphi \circ f$  ( $f \in C_b(X) \otimes E$ ). We show that  $T_\varphi$  is continuous and onto. Let  $G = \{g \in C_b(X) : \|g\| \leq 1\}$  be a neighbourhood of 0 in  $(C_b(X), \|\cdot\|)$ . By continuity of  $\varphi$ , there exists a balanced neighbourhood  $W$  of 0 in  $E$  such that  $|\varphi(a)| \leq 1$  for all  $a \in W$ . Let  $F = \{f \in C_b(X) \otimes E : f(X) \subseteq W\}$ . Then  $T_\varphi(F) \subseteq G$  and so  $T_\varphi$  is continuous. Let  $g \in C_b(X)$ . Choose  $c \in E$  with  $\varphi(c) = r \neq 0$ . Then  $g \otimes \frac{1}{r}c \in C_b(X) \otimes E$  and  $T_\varphi\left(g \otimes \frac{1}{r}c\right) = g$ . Thus  $(C_b(X), \|\cdot\|)$  is

1980 *Mathematics Subject Classification*. Primary 46E10; Secondary 46E40.

*Key words and phrases*. Space of continuous vector-valued functions, uniform topology, separable space,  $\sigma$ -compact metric space, finite covering dimension, approximation property.

separable and so, by Theorem 1,  $X$  is a compact metric space. Now, for any fixed  $z \in X$ , define  $S_z: (C_b(X) \otimes E, \mathcal{U}) \rightarrow E$  by  $S_z(f) = f(z)$ . Then  $S_z$  is also continuous and onto; consequently,  $E$  is separable.

Conversely, suppose  $X$  is a compact metric space and  $E$  is a separable TVS. By Theorem 1,  $(C_b(X), \|\cdot\|)$  is separable. Let  $A = \{g_p\}$  and  $B = \{b_q\}$  be countable dense subsets of  $(C_b(X), \|\cdot\|)$  and  $E$ , respectively. Let  $H$  be the countable subspace generated by  $\{g_p \otimes b_q: p, q = 1, 2, \dots\}$  over rationals. We show that  $H$  is  $\mathcal{U}$ -dense in  $C_b(X) \otimes E$ . Let  $f = \sum_{i=1}^n f_i \otimes a_i$  ( $f_i \in C_b(X)$ ,  $a_i \in E$ ) be in  $C_b(X) \otimes E$  and  $U$  a neighbourhood of 0 in  $E$ . There exists a balanced neighbourhood  $V$  of 0 with  $V + V + \dots + V$  ( $2n$ -terms)  $\subseteq U$ . Choose  $\lambda \geq 1$  such that each  $a_i \in \lambda V$ . Let  $\mu = \max \{\|f_i\|: 1 \leq i \leq n\}$ . For each  $i = 1, \dots, n$ , we can choose  $g_{p_i} \in A$  and  $b_{q_i} \in B$  such that

$$\|g_{p_i} - f_i\| < \frac{1}{\lambda(\mu+1)} \quad \text{and} \quad b_{q_i} - a_i \in \frac{1}{\lambda(\mu+1)} V.$$

Let  $g = \sum_{i=1}^n g_{p_i} \otimes b_{q_i}$ . Then  $g \in H$  and, for any  $x \in X$ ,

$$\begin{aligned} g(x) - f(x) &= \sum_{i=1}^n g_{p_i}(x)(b_{q_i} - a_i) + \sum_{i=1}^n (g_{p_i}(x) - f_i(x))a_i \in \\ &\in \frac{1}{\lambda(\mu+1)} \sum_{i=1}^n g_{p_i}(x)V + \sum_{i=1}^n (g_{p_i}(x) - f_i(x))V \subseteq \\ &\subseteq \frac{1}{\lambda}(V + \dots + V \text{ (} n \text{-terms)}) + \frac{1}{(\mu+1)}(V + \dots + V \text{ (} n \text{-terms)}) \subseteq U. \end{aligned}$$

Hence  $(C_b(X) \otimes E, \mathcal{U})$  is separable. This completes the proof.

Next, using Theorem 2 instead of Theorem 1, we can similarly prove

**THEOREM 4.** *Let  $X$  be locally compact. Then  $(C_0(X) \otimes E, \mathcal{U})$  is separable iff  $X$  is a  $\sigma$ -compact metric space and  $E$  is separable.*

**COROLLARY 5.** *Suppose that either  $X$  has finite covering dimension or  $E$  is locally convex, or  $E$  has the approximation property. Then*

- (i)  $(C_b(X, E), \mathcal{U})$  is separable iff  $X$  is a compact metric space and  $E$  is separable.
- (ii) If  $X$  is locally compact,  $(C_0(X, E), \mathcal{U})$  is separable iff  $X$  is a  $\sigma$ -compact metric space and  $E$  is separable.

**PROOF.** By [2], each of the above restrictions on  $X$  or  $E$  implies that, in case (i),  $C_b(X) \otimes E$  is  $\mathcal{U}$ -dense in  $C_b(X, E)$  and that, in case (ii),  $C_0(X) \otimes E$  is  $\mathcal{U}$ -dense in  $C_0(X, E)$ . The result now follows by applying Theorems 3 and 4.

#### REFERENCES

- [1] SEMADENI, Z. and ZBIJEWSKI, P., Spaces of continuous functions I, *Studia Math.* 16 (1957), 130–141. MR 19–1182.
- [2] SHUCHAT, A. H., Approximation of vector-valued continuous functions, *Proc. Amer. Math. Soc.* 31 (1972), 97–103. MR 44#7267.

- [3] SUMMERS, W. H., Separability in the strict and substrict topologies, *Proc. Amer. Math. Soc.* 35 (1972), 507—514. *MR* 46#2410.
- [4] WARNER, S., The topology of compact convergence on continuous function spaces, *Duke Math. J.* 25 (1958), 265—282. *MR* 21#1521.

(Received December 27, 1983)

DEPARTMENT OF MATHEMATICS  
FACULTY OF SCIENCE  
GARYOUNIS UNIVERSITY  
P.O. BOX 9480  
BENGHAZI  
LIBYA

Current address:

DEPARTMENT OF MATHEMATICS  
FEDERAL GOVERNMENT COLLEGE  
H—8 ISLAMABAD  
PAKISTAN





# PUTTING CONVERGENT SEQUENCES INTO MEASURABLE SETS

S. J. EIGEN

Let  $x_1 > x_2 > \dots > x_n$  be a finite sequence of positive numbers. It is well-known that every measurable subset of the unit interval with positive measure contains a finite sequence similar (i.e. a linear contraction) to the given sequence. This is easily proved using the Lebesgue Density Theorem. P. Erdős [2] raised the question as to what happens if we replace the finite sequence by an infinite sequence which decreases to 0. This is part of an older, and more challenging problem of P. Erdős [3]: given any sequence of positive numbers converging to 0, there is a set of positive measure not containing a sequence similar to the first. In this note, we will show that if the sequence “converges slowly enough” then there exists a set of positive measure which contains no similar sequence. This extends, in part, a result of P. Komjáth [4]: given  $\varepsilon > 0$  there is a set  $H$  of measure  $1 - \varepsilon$  and there is a sequence  $x_n$  converging to 0 such that for each  $t$  in  $[0, 1]$  and every  $a \neq 0$  we have that  $t + ax_n$  is not in  $H$  for infinitely many  $n$ . P. Komjáth begins with a set  $K$  of a certain form, and then constructs the sequence  $x_n$ . The sequence thus constructed is easily seen to converge slowly enough. In this paper, we begin with any sequence which converges slowly enough, and then construct the set  $H$ . Our proof thus works for a larger family of sequences, but we have only that  $t + ax_n$  misses  $H$  for some  $n$ , not necessarily for infinitely many.

Let  $x_n$  be a sequence of positive numbers decreasing to 0. We assume the sequence “converges slowly enough” — by which we mean that the ratios

$$r_n = \frac{x_n}{x_n - x_{n+1}}$$

satisfy  $\lim r_n = \infty$ . The sequence  $1/n$  is an example of a sequence which converges slowly enough. A sequence  $y_n$  is similar to the sequence  $x_n$  if there exist numbers  $a, b, a \neq 0$  such that  $ax_n + b = y_n$  for all  $n$ . Observe that the “tail”  $y_{n+1} - \lim y_i$  to “gap”  $y_n - y_{n-1}$  ratio is  $r_n$ . We will call the absolute value  $|a|$  the “sizing factor”.

LEMMA. Let  $N > 1$  and  $M > 1$  be integers. Then there exists a subset  $A$  of the unit interval with measure  $1 - 1/N$ , and the only sequences in  $A$  similar to the sequence  $x_n$  have a sizing factor less than  $1/M$ .

PROOF. Choose  $n_1$  so that for all  $n \geq n_1$  we have  $r_n > N$ . Choose  $K$  so that  $\frac{1}{K} \leq (x_{n_1})/M$ . Partition the unit interval into the  $K$  pieces  $\left[\frac{i}{K}, \frac{i+1}{K}\right)$   $0 \leq i < K$ . Remove from each the smaller piece  $\left[\frac{(i+1)N-1}{KN}, \frac{i+1}{K}\right)$  of size  $1/KN$ . The

remaining set will be our set  $A$ . Since we removed  $K$  sets of size  $1/KN$ , the set  $A$  has measure  $1 - 1/N$ . In order to see that  $A$  has the second property, observe that  $A$  consists of  $K$  intervals separated by gaps, and the ratio of the interval size to the gap size is  $N - 1$ . Suppose  $y_n$  is a sequence similar to  $x_n$  and contained in  $A$ . If the sequence is wholly contained in a single interval of  $A$ , then by the choice of  $K$  its sizing factor must be less than  $1/M$ . On the other hand, if it is not wholly contained in a single interval then its tail must still be contained in a single interval piece. Thus there must be an  $L$  such that  $y_n$ , for  $n > L$ , are all in the same interval, and  $y_L$  is in a different interval. So, the distance between  $y_L$  and  $y_{L+1}$  is greater than the gap size. Thus the interval length must be greater than

$$y_{L+1} - \lim y_n = (y_L - y_{L+1})r_n.$$

If  $L \equiv n_1$ , then the gap-interval ratio is too large and the tail cannot fit into a single interval. Hence  $L < n_1$ . This means the tail, which must fit into one of the intervals, starts at  $n_1$  or sooner. By the choice of  $K$  this makes the sizing factor less than  $1/M$ .

The existence of the set  $E$  postulated in the beginning follows easily from the lemma.

I would like to thank P. Komjáth for bringing to my attention the references below.

#### REFERENCES

- [1] BORWEIN, D. and DITOR, S. Z., Translates of sequences in sets of positive measure, *Canad. Math. Bull.* **21** (1978), 497—498. MR 80i: 28018.
- [2] ERDŐS, P., Some measure-theoretic problems of a combinatoric and geometric nature, Measure Theory Conference at Oberwolfach, 1983.
- [3] ERDŐS, P., Set theoretic, measure theoretic, combinatorial, and number theoretic problems concerning point sets in Euclidean space, *Real Anal. Exchange* **4** (1978—79), 113—138. MR 80g: 04005.
- [4] KOMJÁTH, P., Large sets not containing images of a given sequence, *Canadian Math. Bull.* **26** (1983), 41—43.

(Received January 19, 1984)

DEPARTMENT OF MATHEMATICS  
COLLEGE OF ARTS AND SCIENCES  
NORTHEASTERN UNIVERSITY  
360 HUNTINGTON AVENUE  
BOSTON, MA 02115  
U.S.A.

# ON APPROXIMATION OF SOLUTIONS OF EXTERIOR BOUNDARY VALUE PROBLEMS

L. SIMON

In [1] and [2] it has been proved that the solution to Dirichlet problem in the exterior of a bounded domain  $D$  for the equation

$$(\Delta + k^2)u = f$$

with homogeneous boundary condition can be obtained as the limit (as  $N \rightarrow \infty$ ) of the solutions of problems in  $\mathbb{R}^3$  for the equations

$$(\Delta + k^2 - V)u = f,$$

where

$$V = \begin{cases} N & \text{in } D, \\ 0 & \text{in } \mathbb{R}^3 \setminus D. \end{cases}$$

A similar theorem has been proved in [3] for the same equation in the case of Neumann problem.

The aim of this paper is to establish analogous results for  $2m$  order linear elliptic equations (with smooth coefficients which are constant out of a bounded set) with general boundary conditions.

## § 1. Preliminaries

For any domain  $G_0 \subset \mathbb{R}^n$  and any nonnegative integer  $k$  denote by  $H^k(G_0)$  the Sobolev space of (complex valued) functions  $f$  whose distributional derivatives of order  $\leq k$  belong to  $L^2(G_0)$ . The norm in  $H^k(G_0)$  is defined by

$$\|f\|_{H^k(G_0)} = \left\{ \sum_{|\alpha| \leq k} \int_{G_0} |D^\alpha f|^2 dx \right\}^{1/2}$$

where  $D = \left( -i \frac{\partial}{\partial x_1}, \dots, -i \frac{\partial}{\partial x_n} \right)$  and  $\alpha = (\alpha_1, \dots, \alpha_n)$  is a multiindex. The definition of  $H^s(G_0)$  for any real  $s$  can be found in [4]. By  $H_{loc}^k(G_0)$  will be denoted the set of functions  $f$  such that  $\varphi f \in H^k(G_0)$  for any  $\varphi \in C_0^\infty(G_0)$ , i.e. for any infinitely differentiable function  $\varphi$  with compact support, contained in  $G_0$ .

Further by  $\Omega$  will be denoted an unbounded domain in  $\mathbb{R}^n$  with bounded smooth boundary  $\partial\Omega$ . Suppose that  $B_{R_0} \supset \bar{G}$  where  $B_{R_0} = \{x \in \mathbb{R}^n : |x| < R_0\}$  and  $G = \mathbb{R}^n \setminus \bar{\Omega}$ .

1980 *Mathematics Subject Classification*. Primary 35J40.

*Key words and phrases*. Exterior boundary value problems for elliptic equations of higher order, interface problems.

Denote by  $H_{loc}^k(\bar{\Omega})$  the set of functions  $f$  with the following property:  $f \in H^k(\Omega \cap B_R)$  for any  $R > R_0$ . For any  $s \geq 0$  we shall use the usual notation  $H^s(\partial\Omega)$  of the Sobolev space of functions defined on  $\partial\Omega$  (see e.g. [4]).

Let  $H^k(\mathbf{R}^n, \partial\Omega)$  be the direct sum of  $H^k(\Omega)$  and  $H^k(\mathbf{R}^n \setminus \bar{\Omega})$ , i.e. the closure in  $\|\cdot\|_{H^k(\mathbf{R}^n)}$  of the linear space of functions which are infinitely differentiable in  $\bar{\Omega}$  and in  $\mathbf{R}^n \setminus \bar{\Omega}$ . (See [5].) Finally, denote by  $H_{loc}^k(\mathbf{R}^n, \partial\Omega)$  the set of functions  $f$  with the property:  $\varphi f \in H^k(\mathbf{R}^n, \partial\Omega)$  for any  $\varphi \in C_0^\infty(\mathbf{R}^n)$ .

Consider an elliptic differential operator  $A$  of order  $2m$  defined by the formula

$$(1.1) \quad Au = \sum_{|\alpha|, |\beta| \leq m} D^\alpha (a_{\alpha\beta} D^\beta u)$$

where  $a_{\alpha\beta}$  are infinitely differentiable (complex valued) functions in  $\mathbf{R}^n$  such that for  $|x| \geq a$   $a_{\alpha\beta}(x)$  are constants:  $a_{\alpha\beta}(x) = a_{\alpha\beta}^0$  ( $a \geq R_0$ ). Thus  $A$  is uniformly elliptic in  $\mathbf{R}^n$ . Let

$$(1.2) \quad P(\xi) = \sum_{|\alpha|, |\beta| \leq m} a_{\alpha\beta}^0 \xi^{\alpha+\beta}, \quad \xi \in \mathbf{R}^n,$$

then  $P(D)$  is a differential operator with constant coefficients which is equal to  $A$  out of  $B_a$ .

Suppose that  $P$  satisfies the assumptions of [8], i.e.

- I.  $S = \{\xi \in \mathbf{R}^n : P(\xi) = 0\}$  is an  $(n-1)$ -dimensional manifold,  
 $\text{grad } P(\xi) \neq 0$  for  $\xi \in S$  and  
 in every point of  $S$  the complete deviation is not 0.

Assumptions I imply that  $S$  consists of a finite number of smooth closed convex  $(n-1)$ -dimensional manifolds  $K_1, \dots, K_r$  (see [8]). Define a direction of the normal  $v$  on each  $K_j$ . For a given unit vector  $\omega$  denote by  $\sigma_j(\omega)$  the (unique) point of  $K_j$  in which the normal  $v$  is parallel to  $\omega$  and it has the same direction as  $\omega$ . Further define  $\mu_j(\omega)$  by the scalar product  $\langle \sigma_j(\omega), \omega \rangle$ . Denote by  $W$  the set of functions  $u$  of the form

$$(1.3) \quad u = \sum_{j=1}^r u_j$$

where  $u_j$  satisfy the inequalities

$$(1.4) \quad |u_j(x)| < c|x|^{(1-n)/2}, \quad \left| \frac{\partial u_j(x)}{\partial r} + i\mu_j(\omega)u_j(x) \right| < c|x|^{-n/2}$$

for sufficiently large  $r = |x|$ ;  $\omega = \frac{x}{r}$ .

In [8] it is proved that for any  $f \in L_a^2(\mathbf{R}^n)$  (i.e. for  $f \in L^2(\mathbf{R}^n)$  such that  $f(x) = 0$  a.e. if  $|x| > a$ ) the equation

$$(1.5) \quad P(D)u = f$$

has a unique solution  $u$  in  $H_{loc}^{2m}(\mathbf{R}^n) \cap W$  and  $u = f * E$  where  $E$  is a fundamental solution of  $P(D)$ , belonging to  $W$ .

Let  $B_j$  be differential operators on  $\partial\Omega$  of order  $r_j \leq 2m-1$  with infinitely differentiable coefficients ( $j = 1, \dots, m$ ). Consider the exterior boundary value problem

for  $u \in H_{loc}^{2m}(\bar{\Omega})$

$$(1.6) \quad Au = f \quad \text{in } \Omega,$$

$$(1.7) \quad B_j u|_{\partial\Omega} = g_j, \quad j = 1, \dots, m,$$

$$(1.8) \quad u \in W.$$

Assume that  $A$  and  $B_j$  satisfy the conditions of regular elliptic boundary value problems in [4],  $f \in L_a^2(\Omega)$  (i.e.  $f \in L^2(\Omega)$  and  $f(x)=0$  a.e. for  $|x|>a$ ) and  $g_j \in H^{2m-r_j-1/2}(\partial\Omega)$ .

Further let  $C_j$  be differential operators on  $\partial G = \partial\Omega$  of order  $l_j \leq 2m-1$  with infinitely differentiable coefficients ( $j=1, \dots, m$ ). Consider the boundary value problem in  $G$

$$(1.9) \quad Au = f_1 \quad \text{in } G$$

$$(1.10) \quad C_j u|_{\partial G} = h_j, \quad j = 1, \dots, m.$$

Suppose that it is a regular elliptic boundary value problem,  $f_1 \in L^2(G)$  and  $h_j \in H^{2m-l_j-1/2}(\partial G)$ . Consider the solutions  $u$  of (1.9), (1.10) in  $H^{2m}(G)$ . Set

$$X_1 = H^{2m-l_1-1/2}(\partial G) \times \dots \times H^{2m-l_m-1/2}(\partial G).$$

LEMMA 1. The kernel  $N_1$  of (1.9), (1.10) is a finite dimensional subspace of  $L^2(G)$ .

$$Y_1 = \{(C_1 u, \dots, C_m u)|_{\partial G}: u \in H^{2m}(G), Au = 0 \text{ in } G\}$$

is a closed subspace of  $X_1$ , consisting of those elements of  $X_1$  which satisfy a finite number of orthogonality conditions. Further, for  $h \in Y_1$ ,  $f_1=0$  there is a unique solution  $u = F_1(h) \in H^{2m}(G)$  of (1.9), (1.10) in

$$N_1^\perp = \{v \in L^2(G): v \perp N_1 \text{ in } L^2(G)\}$$

and  $F_1: Y_1 \rightarrow H^{2m}(G)$  is a bounded linear operator. (See [4].)

For a fixed  $\varepsilon > 0$  denote by  $H^k(\Omega, (1+|x|)^{-1-\varepsilon})$  the Sobolev space with the norm

$$\|v\| = \left\{ \sum_{|\alpha| \leq k} \int_{\Omega} |D^\alpha v|^2 (1+|x|)^{-1-\varepsilon} dx \right\}^{1/2}.$$

LEMMA 2. Suppose that  $u_j \in H_{loc}^{2m}(\bar{\Omega}) \cap W$  satisfies the equation  $Au_j = f$  where  $f \in L_a^2(\Omega)$  and there exists  $u^* \in H_{loc}^{2m-1/2}(\Omega)$  such that

$$(1.11) \quad \lim_{j \rightarrow \infty} \|u_j - u^*\|_{H^{2m-1/2}(\Omega \cap B_R)} = 0 \quad \text{for any } R > R_0.$$

Then for  $|x| > r > a$

$$(1.12) \quad u^*(x) = \sum_{|\alpha|, |\beta| \leq 2m-1} \int_{S_r} d_{\alpha\beta} D^\alpha u^*(y) D^\beta E(x-y) d\sigma_y,$$

where  $S_r = \{x \in \mathbb{R}^n: |x|=r\}$  and  $d_{\alpha\beta}$  depend only on  $\frac{y}{|y|}$ . Moreover,  $u^* \in W$  and for any fixed number  $\varepsilon > 0$

$$(1.13) \quad \lim_{j \rightarrow \infty} \|u_j - u^*\|_{H^{2m-1}(\Omega, (1+|x|)^{-1-\varepsilon})} = 0.$$

PROOF. By Green's formula

$$(1.14) \quad \begin{aligned} u_j(x) = & \sum_{|\alpha|, |\beta| \leq 2m-1} \int_{S_r} c_{\alpha\beta} D^\alpha u_j(y) D^\beta E(x-y) d\sigma_y + \\ & + \sum_{|\alpha|, |\beta| \leq 2m-1} \int_{S_r} d_{\alpha\beta} D^\alpha u_j(y) D^\beta E(x-y) d\sigma_y \end{aligned}$$

if  $R > |x| > r > a$  where  $c_{\alpha\beta}, d_{\alpha\beta}$  depend only on  $y/|y|$ .  $u_j \in W$ ,  $E \in W$  imply that the first term in the right of (1.14) converges to 0 as  $R \rightarrow \infty$ . (See the proof of Theorem 1.7 of [8].) Thus

$$(1.15) \quad u_j(x) = \sum_{|\alpha|, |\beta| \leq 2m-1} \int_{S_r} d_{\alpha\beta} D^\alpha u_j(y) D^\beta E(x-y) d\sigma_y \quad \text{if } |x| > r.$$

Equality (1.11) implies that

$$(1.16) \quad \lim_{j \rightarrow \infty} \|u_j - u^*\|_{H^{2m-1}(S_r)} = 0$$

whence by use of (1.15) one finds equality

$$(1.17) \quad u^*(x) = \sum_{|\alpha|, |\beta| \leq 2m-1} \int_{S_r} d_{\alpha\beta} D^\alpha u^*(y) D^\beta E(x-y) d\sigma_y.$$

( $u_j$ ) contains a subsequence which converges a.e. to  $u^*$ , consequently, (1.17) holds a.e. and by the continuity of  $u^*$  (1.17) holds in every point.)

From (1.17) it follows that  $u^* \in W$ . Further, from (1.15), (1.17) we obtain for  $|x| > r$  the formulas

$$D^\gamma u_j(x) = \sum_{|\alpha|, |\beta| \leq 2m-1} \int_{S_r} d_{\alpha\beta} D^\alpha u_j(y) D^{\beta+\gamma} E(x-y) d\sigma_y,$$

$$D^\gamma u^*(x) = \sum_{|\alpha|, |\beta| \leq 2m-1} \int_{S_r} d_{\alpha\beta} D^\alpha u^*(y) D^{\beta+\gamma} E(x-y) d\sigma_y,$$

thus

$$(1.18) \quad |D^\gamma u_j(x) - D^\gamma u^*(x)|^2 \leq c_1 \sum_{|\alpha|, |\beta| \leq 2m-1} \left\{ \int_{S_r} |D^\alpha u_j - D^\alpha u^*|^2 d\sigma \right\} \left\{ \int_{S_r} |D^{\beta+\gamma} E(x-y)|^2 d\sigma_y \right\}.$$

Hence

$$(1.19) \quad \begin{aligned} & \int_{\Omega \setminus B_{r+1}} |D^\gamma u_j(x) - D^\gamma u^*(x)|^2 (1 + |x|)^{-1-\varepsilon} dx \leq \\ & \leq c_1 \sum_{|\alpha|, |\beta| \leq 2m-1} \left\{ \int_{S_r} |D^\alpha u_j - D^\alpha u^*|^2 d\sigma \right\} \times \\ & \times \left\{ \int_{\Omega \setminus B_{r+1}} (1 + |x|)^{-1-\varepsilon} \left[ \int_{S_r} |D^{\beta+\gamma} E(x-y)|^2 d\sigma_y \right] dx \right\}. \end{aligned}$$

Since  $E \in W$  thus  $D^{\beta+\gamma} E \in W$  (see [8]) and (for a fixed  $r$ )

$$\int_{S_r} |D^{\beta+\gamma} E(x-y)|^2 d\sigma_y \leq \frac{c_2}{(1 + |x|)^{n-1}} \quad \text{if } |x| > r+1.$$



Consequently, (1.11), (1.16), (1.19) imply (1.13). ■

Let

$$X_2 = H^{2m-r_1-1/2}(\partial\Omega) \times \dots \times H^{2m-r_m-1/2}(\partial\Omega)$$

and denote by  $L^2(\Omega, (1+|x|)^{-1-\varepsilon})$  the  $L^2$  space over  $\Omega$  with the weight function  $x \mapsto (1+|x|)^{-1-\varepsilon}$ . Further, denote by  $\psi$  a fixed function with the properties:

$$(1.20) \quad \psi \in C^\infty(\Omega), \quad \psi > 0, \quad \psi(x) = \begin{cases} 1 & \text{for } |x| \leq a, \\ \exp(-|x|) & \text{for } |x| > 2a. \end{cases}$$

LEMMA 3. The kernel  $N_2$  of (1.6)—(1.8) is a finite dimensional subspace of  $L^2(\Omega, (1+|x|)^{-1-\varepsilon})$ . There exists a closed subspace  $Y_2$  of  $X_2$  (consisting of those elements of  $X_2$  which satisfy a finite number of orthogonality conditions) such that the problem (1.6)—(1.8) with  $f=0$  has a solution  $u \in H_{loc}^{2m}(\bar{\Omega})$  if and only if  $g=(g_1, \dots, g_m) \in Y_2$ . Further, for  $g \in Y_2, f=0$  there is a unique solution  $u=F_2(g)$  of (1.6)—(1.8) in

$$N_2^\perp = \{v \in L^2(\Omega, (1+|x|)^{-1-\varepsilon}), v \perp N_2 \text{ in } L^2(\Omega, (1+|x|)^{-1-\varepsilon})\};$$

$F_2: Y_2 \rightarrow H_{loc}^{2m}(\bar{\Omega})$  is a linear operator and the following estimation holds;

$$(1.21) \quad \|\psi u\|_{H^{2m}(\Omega)} \leq c_1 \sum_{j=1}^m \|g_j\|_{H^{2m-r_j-1/2}(\partial\Omega)}$$

where the constant  $c_1 > 0$  does not depend on  $g$ .

PROOF. According to Theorem 1 of [9] the kernel  $N_2$  of (1.6)—(1.8) is a finite dimensional subspace of  $L^2(\Omega, (1+|x|)^{-1-\varepsilon})$  (if  $Au=0, B_j u|_{\partial\Omega}=0$  then  $u$  is a smooth function and thus  $u \in W$  implies  $u \in L^2(\Omega, (1+|x|)^{-1-\varepsilon})$ ) and the problem (1.6)—(1.8) with  $f=0$  has a solution if  $g=(g_1, \dots, g_m)$  satisfies a finite number of orthogonality conditions. Further, for the solutions  $u$  of (1.6)—(1.8) (with  $f=0$ ) the following apriori estimation holds:

$$(1.22) \quad \|\psi u\|_{H^{2m}(\Omega)} \leq c_1 \left[ \sum_{j=1}^m \|g_j\|_{H^{2m-r_j-1/2}(\partial\Omega)} + \|u\|_{H^{2m-1}(\Omega \cap B_{R_0})} \right]$$

where  $\psi$  is a fixed function, satisfying (1.20) and the constant  $c_1 > 0$  does not depend on  $u$ .

The problem (1.6)—(1.8) may have at most one solution in  $N_2^\perp$ , because if  $u$  satisfies (1.6)—(1.8) with  $f=0, g=0$  and  $u \in N_2^\perp$  then by  $u \in N_2, u \perp u$  in  $L^2(\Omega, (1+|x|)^{-1-\varepsilon})$  ( $u \in W$  implies that  $u \in L^2(\Omega, (1+|x|)^{-1-\varepsilon})$ ) and thus  $u=0$ .

Denote by  $Y_2$  the set of  $g \in X_2$  such that the problem (1.6)—(1.8) with  $f=0$  has a solution. Clearly,  $Y_2$  is a linear subspace of  $X_2$ . For any  $g \in Y_2$  there exists a solution  $u_1 \in N_2^\perp$  of (1.6)—(1.8). Indeed, let  $u$  be a solution of (1.6)—(1.8). Then  $u \in L^2(\Omega, (1+|x|)^{-1-\varepsilon})$  can be written in the form

$$u = u_1 + u_2 \quad \text{where } u_1 \in N_2^\perp \quad \text{and } u_2 \in N_2.$$

Since  $u_2 \in N_2$  thus  $u_1$  is a solution of (1.6)—(1.8), too.

For any  $g \in Y_2$  denote by  $F_2(g)$  the unique solution  $u$  of (1.6)—(1.8) (with  $f=0$ ) in  $N_2^\perp$ . Clearly,  $F_2: Y_2 \rightarrow H_{loc}^{2m}(\bar{\Omega})$  is a linear operator.



To complete the proof we have to show that  $Y_2$  is a closed subspace of  $X_2$  and the estimation (1.21) holds. In order to prove (1.21) it is sufficient to show inequality

$$(1.23) \quad \|u\|_{H^{2m-1}(\Omega \cap B_{R_0})} \leq c_2 \sum_{j=1}^m \|g_j\|_{H^{2m-r_j-1/2}(\partial\Omega)}$$

(see (1.22)). Suppose that (1.23) is not true. Then there exists a sequence of  $u_k \in H_{loc}^{2m}(\bar{\Omega}) \cap W$  such that  $Au_k = 0$ ,  $u_k \in N_2^\perp$  and

$$(1.24) \quad \|u_k\|_{H^{2m-1}(\Omega \cap B_{R_0})} = 1,$$

$$(1.25) \quad \lim_{k \rightarrow \infty} \|B_j u_k\|_{H^{2m-r_j-1/2}(\partial\Omega)} = 0.$$

(1.22), (1.24) and (1.25) imply that  $\|\psi u_k\|_{H^m(\Omega)}$  is bounded.

Consequently, for any fixed  $R > R_0$  there is a subsequence  $(u'_k)$  of  $(u_k)$  such that  $(u'_k)$  converges to a function  $u^* \in H^{2m-1/2}(\Omega \cap B_R)$  with respect to the norm of  $H^{2m-1/2}(\Omega \cap B_R)$ . Applying this argument to  $R = R_0 + 1, R_0 + 2, R_0 + 3, \dots$  we find subsequences  $(u_k^{(1)}), (u_k^{(2)}), (u_k^{(3)}), \dots$  of  $(u_k)$  such that  $(u_k^{(r+1)})$  is a subsequence of  $(u_k^{(r)})$  and  $(u_k^{(r)})$  converges to a function  $u^*$  with respect to the norm of  $H^{2m-1/2}(\Omega \cap B_{R_0+r})$ . Thus  $(u_k^{(k)})$  is a subsequence of  $(u_k)$  converging to a function  $u^* \in H_{loc}^{2m-1/2}(\bar{\Omega})$  with respect to the norm of  $H^{2m-1/2}(\Omega \cap B_R)$  for any  $R > R_0$ .

Thus by making use of  $Au_k^{(k)} = 0$  we find that for any  $R > R_0$

$$\|Au^*\|_{H^{-1/2}(\Omega \cap B_R)} = \lim_{k \rightarrow \infty} \|Au_k^{(k)} - Au^*\|_{H^{-1/2}(\Omega \cap B_R)} = 0$$

(see [4]) and, consequently,  $Au^* = 0$  in  $\Omega$ . Further, (1.25) implies that

$$\|B_j u^*\|_{H^{2m-r_j-1/2}(\partial\Omega)} = \lim_{k \rightarrow \infty} \|B_j u_k^{(k)}\|_{H^{2m-r_j-1/2}(\partial\Omega)} = 0,$$

$B_j u^* = 0$ . Therefore the function  $u^* \in H_{loc}^{2m-1/2}(\bar{\Omega})$  satisfies (1.6) and (1.7) with  $f=0, g=0$ . Since  $A$  and  $B$  satisfy the regularity conditions of [4] thus  $u^* \in H_{loc}^{2m}(\bar{\Omega})$ . Lemma 2 implies that  $u^* \in W$  and so  $u^* \in N_2$ . From Lemma 2 it follows that  $(u_k^{(k)})$  converges to  $u^*$  with respect to the norm of  $L^2(\Omega, (1+|x|)^{-1-\epsilon})$ . As  $u_k^{(k)} \in N_2^\perp$  thus  $u^* \in N_2^\perp$ . Consequently,  $u^* = 0$ , but it is impossible since by (1.24)

$$\|u^*\|_{H^{2m-1}(\Omega \cap B_{R_0})} = \lim_{k \rightarrow \infty} \|u_k^{(k)}\|_{H^{2m-1}(\Omega \cap B_{R_0})} = 1.$$

Thus inequalities (1.23) and (1.21) are proved.

Finally, we show that  $Y_2$  is a closed subspace of  $X_2$ . If a sequence of  $g^{(k)} = (g_1^{(k)}, \dots, g_m^{(k)}) \in Y_2$  converges to  $g = (g_1, \dots, g_m) \in X_2$  with respect to the norm of  $X_2$  then by (1.21)  $F_2(g^{(k)})$  converges to a function  $u \in H_{loc}^{2m}(\bar{\Omega})$  in  $H^{2m}(\Omega \cap B_R)$  for any  $R > R_0$ . Thus  $Au = 0$  and Lemma 2 implies that  $u \in W$ . Further,

$$\|B_j u - g_j\|_{H^{2m-r_j-1/2}(\partial\Omega)} = \lim_{k \rightarrow \infty} \|B_j(F(g^{(k)})) - g_j^{(k)}\|_{H^{2m-r_j-1/2}(\partial\Omega)} = 0,$$

i.e.  $B_j u = g_j$  ( $j=1, \dots, m$ ). Therefore,  $g \in Y_2$  and the proof of Lemma 3 is complete. ■

## § 2. Approximation theorems

Let  $(\lambda_N)$  be a sequence of positive numbers such that  $\lim_{N \rightarrow \infty} \lambda_N = +\infty$ . Consider the problem

$$(2.1) \quad Au_N = f \quad \text{in } \Omega, \quad Au_N = f_1 \quad \text{in } G$$

$$(2.2) \quad (B_j u_N)^+ = \lambda_N (B_j u_N)^-, \quad j = 1, \dots, m$$

$$(2.3) \quad (C_j u_N)^+ = (C_j u_N)^-, \quad j = 1, \dots, m$$

$$(2.4) \quad u_N \in W$$

where by  $g^+$  and  $g^-$  is denoted the trace on  $\partial\Omega = \partial G$  of  $g|_G$  resp.  $g|_\Omega$  and  $A, B_j, C_j, f, f_1$  satisfy the conditions, formulated in § 1.

**THEOREM 1.** Suppose that for any  $f \in L^2_a(\Omega)$ ,  $g_j \in H^{2m-r_j-1/2}(\partial\Omega)$  there exists at least one solution  $u \in H^{2m}_{loc}(\Omega)$  of (1.6)–(1.8) and for any  $f_1 \in L^2(G)$ ,  $h_j \in H^{2m-l_j-1/2}(\partial G)$  there exists at least one solution  $u \in H^{2m}(G)$  of (1.9), (1.10).

Then for sufficiently large  $N$  the problem (2.1)–(2.4) has a unique solution  $u_N \in H^{2m}_{loc}(\mathbb{R}^n, \partial\Omega)$  such that  $u_N|_G \in N_1^\perp$  and  $u_N|_\Omega \in N_2^\perp$ . For arbitrary fixed  $\varepsilon > 0$

$$(2.5) \quad \|u_N - u_0\|_{H^{2m}(\Omega, (1+|x|)^{-1-\varepsilon})} = O\left(\frac{1}{\lambda_N}\right),$$

$$(2.6) \quad \|u_N - u_0\|_{H^{2m}(G)} = O\left(\frac{1}{\lambda_N}\right)$$

where  $u_0|_\Omega \in H^{2m}_{loc}(\Omega)$  is the unique solution of (1.6)–(1.8) with  $g_j = 0$ , belonging to  $N_2^\perp$  and  $u_0|_G \in H^{2m}(G)$  is the unique solution of (1.9), (1.10) with  $h_j = (C_j u_0)^-$ , belonging to  $N_1^\perp$ .

**PROOF.** Assumptions of Theorem 1 imply that (for arbitrary  $f \in L^2_a(\Omega)$ ,  $f_1 \in L^2(G)$ ) there is a unique function  $u_0 \in H^{2m}_{loc}(\mathbb{R}^n, \partial\Omega)$  such that  $u_0|_\Omega$  is the solution of (1.6)–(1.8) with  $g_j = 0$ , belonging to  $N_2^\perp$  and  $u_0|_G$  is the solution of (1.9), (1.10) with  $h_j = (C_j u_0)^-$ , belonging to  $N_1^\perp$ . (See the proof of Lemma 3.) Thus  $u_N \in H^{2m}_{loc}(\mathbb{R}^n, \partial\Omega)$  is a solution of (2.1)–(2.4) — satisfying  $u_N|_G \in N_1^\perp$  and  $u_N|_\Omega \in N_2^\perp$  — if and only if  $v_N = u_N - u_0 \in H^{2m}_{loc}(\mathbb{R}^n, \partial\Omega)$  satisfies

$$(2.7) \quad v_N|_G \in N_1^\perp, \quad v_N|_\Omega \in N_2^\perp,$$

$$(2.8) \quad Av_N = 0 \quad \text{in } \Omega, \quad Av_N = 0 \quad \text{in } G,$$

$$(2.9) \quad (B_j v_N)^+ = \lambda_N (B_j v_N)^- - \varphi_j, \quad j = 1, \dots, m,$$

$$(2.10) \quad (C_j v_N)^+ = (C_j v_N)^-, \quad j = 1, \dots, m,$$

$$(2.11) \quad v_N \in W$$

where  $\varphi_j$  is defined by  $\varphi_j = (B_j u_0)^+$ .

If  $v_N \in H_{loc}^{2m}(\mathbf{R}^n, \partial\Omega)$  is a solution of (2.7)–(2.11) then  $w_N = ((B_1 v_N)^-, \dots, (B_m v_N)^-) \in X_2$  satisfies

$$(2.12) \quad w_N - \frac{1}{\lambda_N} L w_N = \frac{1}{\lambda_N} \varphi$$

where  $\varphi = (\varphi_1, \dots, \varphi_m) \in X_2$ ,  $L = B F_1 C F_2$ ;  $F_1, F_2$  are defined in Lemma 1 resp. in Lemma 3 and  $B, C$  are defined by

$$Bu = ((B_1 u)^+, \dots, (B_m u)^+), \quad Cu = ((C_1 u)^-, \dots, (C_m u)^-).$$

Further, if  $w_N \in X_2$  is a solution of (2.12) then

$$(2.13) \quad v_N = \begin{cases} F_2 w_N & \text{in } \Omega \\ F_1 [C(F_2 w_N)] & \text{in } G \end{cases}$$

satisfies (2.7)–(2.11).

Since  $C$  maps  $H^{2m}(\Omega \cap B_{R_0})$  into  $X_1$  and  $B$  maps  $H^{2m}(G)$  into  $X_2$  linearly and continuously thus by Lemma 1 and Lemma 3.

$$L: X_2 \rightarrow X_2$$

is a continuous linear operator. (By assumption of Theorem 1  $F_1$  and  $F_2$  are defined on the whole  $X_1$  resp.  $X_2$ .) Consequently, for  $\lambda_N > \|L\|$  there exists a unique solution  $w_N$  of (2.12) in  $X_2$  and thus there exists a unique solution  $v_N$  of (2.7)–(2.11). This solution can be given by (2.13).

Further,

$$(2.14) \quad \begin{aligned} \|w_N\|_{X_2} &= \left\| \left( I - \frac{L}{\lambda_N} \right)^{-1} \frac{\varphi}{\lambda_N} \right\|_{X_2} = \left\| \sum_{k=0}^{\infty} \frac{1}{\lambda_N^k} L^k \left( \frac{\varphi}{\lambda_N} \right) \right\|_{X_2} \leq \\ &\leq \frac{1}{\lambda_N} \|\varphi\|_{X_2} \sum_{k=0}^{\infty} \frac{\|L\|^k}{\lambda_N^k} \leq \frac{2}{\lambda_N} \|\varphi\|_{X_2} \end{aligned}$$

for  $\lambda_N > 2\|L\|$ . Consequently, by Lemma 3 and (2.13), for any fixed  $R > R_0$

$$(2.15) \quad \begin{aligned} \|u_N - u_0\|_{H^{2m}(\Omega \cap B_R)} &= \|v_N\|_{H^{2m}(\Omega \cap B_R)} = \\ &= \|F_2 w_N\|_{H^{2m}(\Omega \cap B_R)} \leq c_1(R) \|w_N\|_{X_2} \leq \frac{2c_1(R)}{\lambda_N} \|\varphi\|_{X_2} \end{aligned}$$

(where  $c_1(R)$  does not depend on  $\lambda_N$  and  $\varphi$ ). Thus estimation (1.19) implies (2.5).

Finally, from (2.13), (2.14), Lemma 1 and Lemma 3 we obtain that

$$(2.16) \quad \|u_N - u_0\|_{H^{2m}(G)} = \|v_N\|_{H^{2m}(G)} = \|F_1 [C(F_2 w_N)]\|_{H^{2m}(G)} \leq \frac{c_2}{\lambda_N} \|\varphi\|_{X_2}$$

(the number  $c_2 > 0$  does not depend on  $\lambda_N$  and  $\varphi$ ) whence one obtains (2.6). ■

REMARK 1. By making use of Theorem 1 of [9] it can be proved that for any fixed  $R > 0$

$$\|u_0\|_{H^{2m}(\Omega \cap B_R)} \leq c_2(R) \|f\|_{L^2_+(\Omega)},$$

whence

$$\|u_0\|_{H^{2m}(G)} \leq c_4 [\|f\|_{L^2_a(\Omega)} + \|f_1\|_{L^2(G)}]$$

and so

$$\|\varphi\|_{X_2} \leq c_5 [\|f\|_{L^2_a(\Omega)} + \|f_1\|_{L^2(G)}].$$

Consequently, from (1.19), (2.15), (2.16) we find the estimations

$$\|u_N - u_0\|_{H^{2m}(\Omega, (1+|x|)^{-1-\nu})} \leq \frac{c_6}{\lambda_N} [\|f\|_{L^2_a(\Omega)} + \|f_1\|_{L^2(G)}],$$

$$\|u_N - u_0\|_{H^{2m}(G)} \leq \frac{c_7}{\lambda_N} [\|f\|_{L^2_a(\Omega)} + \|f_1\|_{L^2(G)}].$$

**THEOREM 2.** Let  $f \in L^2_a(\Omega)$  and  $f_1 \in L^2(G)$  be functions such that there exists  $\tilde{u} \in W \cap H^{2m}_{loc}(\mathbb{R}^n, \partial\Omega)$  with the property

$$A\tilde{u} = f \quad \text{in } \Omega, \quad A\tilde{u} = f_1 \quad \text{in } G.$$

Further, let  $(\lambda_N)$  be a sequence of positive numbers such that  $\lim_{N \rightarrow \infty} \lambda_N = +\infty$  and for all  $N$  there exists a solution  $u_N \in H^{2m}_{loc}(\mathbb{R}^n, \partial\Omega)$  of (2.1)–(2.4) with the property  $u_N|_G \in N_1^\perp$  and  $u_N|_\Omega \in N_2^\perp$ .

Then there exists a (unique) function  $u_0 \in H^{2m}_{loc}(\mathbb{R}^n, \partial\Omega)$  such that  $u_0|_\Omega$  satisfies (1.6)–(1.8) with  $g_j = 0$ ,  $u_0|_\Omega \in N_2^\perp$  and  $u_0|_G$  satisfies (1.9), (1.10) with  $h_j = (C_j u_0)^-$ ,  $u_0|_G \in N_1^\perp$ . Finally, for these solutions  $u_N, u_0$  the estimations (2.5), (2.6) hold.

**PROOF.** Let  $\tilde{u} = \tilde{u}_0 + \tilde{u}_1$ , where

$$\tilde{u}_0|_\Omega \in N_2^\perp, \quad \tilde{u}_1|_\Omega \in N_2; \quad \tilde{u}_0|_G \in N_1^\perp, \quad \tilde{u}_1|_G \in N_1.$$

Then

$$A\tilde{u}_0 = f \quad \text{in } \Omega, \quad A\tilde{u}_0 = f_1 \quad \text{in } G$$

and, consequently,  $v_N = u_N - \tilde{u}_0 \in H^{2m}_{loc}(\mathbb{R}^n, \partial\Omega)$  satisfies

$$(2.17) \quad v_N|_G \in N_1^\perp, \quad v_N|_\Omega \in N_2^\perp,$$

$$(2.18) \quad Av_N = 0 \quad \text{in } \Omega, \quad Av_N = 0 \quad \text{in } G,$$

$$(2.19) \quad (B_j v_N)^+ = \lambda_N (B_j v_N)^- + \lambda_N p_j - q_j, \quad j = 1, \dots, m,$$

$$(2.20) \quad (C_j v_N)^+ = (C_j v_N)^- + r_j, \quad j = 1, \dots, m,$$

$$(2.21) \quad v_N \in W$$

where

$$p_j = (B_j \tilde{u}_0)^-, \quad q_j = (B_j \tilde{u}_0)^+, \quad r_j = (C_j \tilde{u}_0)^- - (C_j \tilde{u}_0)^+.$$

Define the extensions  $\tilde{F}_1: X_1 \rightarrow H^{2m}(G)$  and  $\tilde{F}_2: X_2 \rightarrow H^{2m}_{loc}(\bar{\Omega})$  of  $F_1$  resp.  $F_2$  by

$$\tilde{F}_1 = F_1 P_1, \quad \tilde{F}_2 = F_2 P_2$$

where by  $P_k$  is denoted the projection from  $X_k$  onto the closed subspace  $Y_k$ . Then by Lemma 1 and Lemma 3  $\tilde{F}_1$  is linear and continuous,  $\tilde{F}_2$  is a linear operator such

that for any  $g \in X_2$  the estimation

$$\|\psi \tilde{F}_2 g\|_{H^2m(\Omega)} \leq c_1 \|g\|_{X_2}$$

holds. Therefore  $L = B\tilde{F}_1 C\tilde{F}_2: X_2 \rightarrow X_2$  is a continuous linear operator. (See the proof of Theorem 1.)

Since  $v_N$  is a solution of (2.17)–(2.21) thus

$$(2.22) \quad w_N = ((B_1 v_N)^-, \dots, (B_m v_N)^-) \in X_2$$

satisfies

$$(2.23) \quad w_N - \frac{1}{\lambda_N} L w_N = s_N$$

where

$$s_N = \frac{1}{\lambda_N} q - p + \frac{1}{\lambda_N} B\tilde{F}_1 r,$$

$$p = (p_1, \dots, p_m), \quad q = (q_1, \dots, q_m), \quad r = (r_1, \dots, r_m).$$

Hence for  $\lambda_N > \|L\|$

$$w_N = \left( I - \frac{1}{\lambda_N} L \right)^{-1} s_N = \sum_{k=0}^{\infty} \frac{1}{\lambda_N^k} L^k s_N$$

and thus for  $\lambda_N > 2\|L\|$

$$\begin{aligned} \|w_N + p\|_{X_2} &\leq \|w_N - s_N\|_{X_2} + \|s_N + p\|_{X_2} \\ (2.24) \quad &\leq \left\| \sum_{k=1}^{\infty} \frac{1}{\lambda_N^k} L^k s_N \right\|_{X_2} + \frac{1}{\lambda_N} \|q + B\tilde{F}_1 r\|_{X_2} \\ &\leq \frac{2}{\lambda_N} (\|L\| \|p\|_{X_2} + \|q\|_{X_2} + \|B\tilde{F}_1 r\|_{X_2}) \end{aligned}$$

because

$$\begin{aligned} \left\| \sum_{k=1}^{\infty} \frac{1}{\lambda_N^k} L^k s_N \right\|_{X_2} &\leq \left( \sum_{k=1}^{\infty} \frac{\|L\|^k}{\lambda_N^k} \right) \|s_N\|_{X_2} \leq \frac{2\|L\|}{\lambda_N} \|s_N\|_{X_2} \\ &\leq \frac{1}{\lambda_N} (2\|L\| \|p\|_{X_2} + \|q\|_{X_2} + \|B\tilde{F}_1 r\|_{X_2}). \end{aligned}$$

Consequently,

$$(2.25) \quad \lim_{N \rightarrow \infty} \|w_N + p\|_{X_2} = 0.$$

Since by (2.17), (2.18), (2.21), (2.22)  $w_N \in Y_2$  (see Lemma 3) and  $Y_2$  is a closed subspace of  $X_2$  thus (2.25) implies that  $-p \in Y_2$ . Further, as by (2.17), (2.18), (2.20)–(2.22)  $C\tilde{F}_2 w_N + r \in Y_1$  (see Lemma 1) and  $Y_1$  is a closed subspace of  $X_1$  thus

$$\lim_{N \rightarrow \infty} \|C\tilde{F}_2 w_N + C\tilde{F}_2 p\|_{X_1} = 0$$

implies that  $C\tilde{F}_2(-p) + r \in Y_1$ . Consequently, by Lemma 1 and Lemma 3

$$(2.26) \quad v_0 = \begin{cases} \tilde{F}_2(-p) & \text{in } \Omega \\ \tilde{F}_1[C(\tilde{F}_2(-p))] + \tilde{F}_1 r & \text{in } G \end{cases}$$

satisfies

$$v_0|_{\Omega} \in N_2^{\perp}, \quad v_0|_G \in N_1^{\perp},$$

$$Av_0 = 0 \quad \text{in } \Omega,$$

and

$$(B_j v_0)^- = -p_j, \quad j = 1, \dots, m, \quad v_0 \in W$$

$$Av_0 = 0 \quad \text{in } G,$$

$$(C_j v_0)^+ = (C_j v_0)^- + r_j, \quad j = 1, \dots, m.$$

Hence, by making use of the definition of  $p_j$  and  $r_j$  we obtain that for  $u_0 = v_0 + \tilde{u}_0 \in H_{loc}^{2m}(\mathbf{R}^n, \partial\Omega)$ ,  $u_0|_{\Omega}$  satisfies (1.6)–(1.8) with  $g_j = 0$ ,  $u_0|_{\Omega} \in N_2^{\perp}$  and  $u_0|_G$  satisfies (1.9), (1.10) with  $h_j = (C_j u_0)^-$ ,  $u_0|_G \in N_1^{\perp}$ .

Since by (2.17), (2.18), (2.20)–(2.22) and (2.26)

$$u_N - u_0 = v_N - v_0 = \begin{cases} \tilde{F}_2(w_N + p) & \text{in } \Omega \\ \tilde{F}_1[C(\tilde{F}_2(w_N + p))] & \text{in } G \end{cases}$$

thus (2.24), Lemma 1 and Lemma 3 imply the estimations (2.5), (2.6). ■

REMARK 2. Similar theorems can be proved in the case if (instead of assumptions I) the polynomial  $P$  satisfies the assumptions in [6] resp. [7], i.e.

II.  $P(\xi) \neq 0$  for  $\xi \in \mathbf{R}^n \setminus \{0\}$ ;

$$P(D) = \sum_{j=1}^{2in} P_j(D)$$

where  $P_j(D)$  denotes the homogeneous part of  $P(D)$  of order  $j$  and

$$P_l(\xi) \neq 0 \quad \text{for } \xi \in \mathbf{R}^n \setminus \{0\}; \quad l_i < \frac{n}{2} + 1.$$

#### REFERENCES

- [1] Рамм, А. Г., Об одном методе решения задачи Дирихле в бесконечной области, *Izv. Vysš. Učebn. Zaved. Matematika* 5 (1965), 124–127. *MR* 32# 7933.
- [2] LAX, P. D. and PHILLIPS, R. S., Decaying modes for the wave equation in the exterior of an obstacle, *Comm. Pure Appl. Math.* 22 (1969), 737–788. *MR* 40# 7641.
- [3] RAMM, A. G., Exterior boundary value problems as limits of interface problems, *J. Math. Anal. Appl.* 84 (1981), 256–263. *MR* 83a: 35025.
- [4] LIONS, J. L. and MAGENES, E., *Problèmes aux limites non homogènes et applications*, vol. 1, Dunod, Paris, 1968. *MR* 40# 512.
- [5] Шефтель, З. Г., Энергетические неравенства и общие граничные задачи для эллиптических уравнений с разрывными коэффициентами, *Sibirsk. Mat. Ž.* 6 (1965), 636–668. *MR* 31# 6056.
- [6] Шимон, Л., Об аппроксимации решений граничных задач в неограниченных областях, *Дифференциальные Уравнения* 9 (1973), 1482–1492. *MR* 52# 6178.

- [7] SIMON, L., On elliptic differential equations in  $R^n$ , *Ann Univ. Sci. Budapest. Eötvös Sect Math.* 27 (1985), 241—256.
- [8] Вайнберг, Б. Р., Принципы излучения, предельного поглощения и предельной амплитуды в общей теории уравнений с частными производными, *Uspehi Mat. Nauk* 21 (1966), 115—194. *MR* 35#4559.
- [9] Вайнберг, Б. Р., Об эллиптических задачах в неограниченных областях, *Mat. Sb. (N. S.)* 75 (117) (1968), 3, 454—480. *MR* 34#6601.

( Received February 6, 1984 )

EÖTVÖS LORÁND TUDOMÁNYEGYETEM  
TERMÉSZETTUDOMÁNYI KAR  
ANALÍZIS TANSZÉK  
MŰZEUM KRT. 6—8  
H—1088 BUDAPEST  
HUNGARY



## A CHARACTERIZATION OF GENERALIZED MATROID LATTICES

P. DAMASCHKE and M. STERN

### 1. Introduction

In [7] generalized matroid lattices have been introduced as lattices  $L$  of finite length which possess the isomorphism property

$$[u \wedge b, u] \cong [b, u \vee b]$$

for all join-irreducible elements  $u \in L$  and for all  $b \in L$ .

In [4] generalized matroid lattices have been characterized as lattices satisfying  $(b, u)M$  and  $(u, b)M^*$  for each join-irreducible element  $u \in L$  and for each element  $b \in L$  (concerning this notation, we refer to Section 2).

It is the aim of the present paper to show that the condition  $(b, u)M$  follows already from  $(u, b)M^*$  and can therefore be dropped.

We prove this by showing first that  $(u, b)M^*$  implies (upper) semimodularity (cf. Theorem 2) and even strong semimodularity (cf. Corollary 4). Finally we show that  $(u, b)M^*$  is equivalent to the abovementioned isomorphism property (cf. Theorem 6).

### 2. Preliminaries

First we need the concept of (dual) modular pairs in lattices:

**DEFINITION.** Let  $L$  be a lattice and  $a, b \in L$ . We say that  $(a, b)$  is a modular pair, and we write  $(a, b)M$  if

$$c \leq b \text{ implies } (c \vee a) \wedge b = c \vee (a \wedge b) \quad (c \in L).$$

If  $(a, b)$  is not a modular pair, we write  $(a, b)\overline{M}$ . We say that  $(a, b)$  is a dual-modular pair, and we write  $(a, b)M^*$  if

$$c \leq b \text{ implies } (c \wedge a) \vee b = c \wedge (a \vee b) \quad (c \in L).$$

If  $(a, b)$  is not a dual-modular pair, we write  $(a, b)\overline{M^*}$ .

Throughout this paper let now  $L$  always be a lattice of finite length.

We write  $x < y$  if  $x$  is a lower cover of  $y$ . An element  $u \in L$  is called join-irreducible if it has exactly one lower cover which will be denoted by  $u'$ .

---

1980 *Mathematics Subject Classification*. Primary 06D99; Secondary 05B35.

*Key words and phrases*. Matroid lattices, semimodular lattices of finite length, join-irreducible elements, dual-modular pairs, isomorphism property for join-irreducible elements.

By  $J(L)$  we denote the set of all join-irreducible elements of  $L$ . An induction argument on the length of  $L$  shows that each element ( $\neq 0$ ) is a join of join-irreducible elements.

The present note deals with the following

QUESTION. What does it mean for  $L$  that

$$(*) \quad u \in J(L) \text{ implies } (u, b)M^* \quad (\forall b \in L)?$$

First we state that from  $(*)$  it follows that  $(p, b)M^*$  holds for each atom  $p \in L$  and for each  $b \in L$ . From this we get by [5, Theorem 7.6, p. 31] that  $L$  has the covering property

$$(C) \quad p \text{ atom, } b \text{ arbitrary element and } b \wedge p = 0 \text{ imply } b < b \vee p.$$

The covering property (C) does not, in general, imply the neighbourhood condition

$$(N) \quad a \wedge b < a \Rightarrow b < a \vee b.$$

This can be seen if we consider, for example, the lattice of Fig. 1:

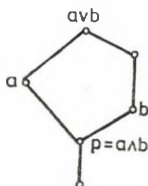


Fig. 1

In this lattice, the covering property (C) trivially holds. On the other hand, we have  $a \wedge b < b$  but  $b$  is not a lower cover of  $a \vee b$ , i.e. the neighbourhood condition (N) does not hold.

We remark that in lattices of finite length, the neighbourhood condition (N) means nothing else than (upper) semimodularity.

We shall see that lattices satisfying  $(*)$  are (strongly) semimodular (cf. Theorem 2 and Corollary 4). More precisely, we shall even show that  $(*)$  is equivalent to the isomorphism property

$$[u \wedge b, u] \cong [b, b \vee u]$$

for each  $u \in J(L)$  and for each  $b \in L$  (cf. Theorem 6). This means that the lattices satisfying  $(*)$  are precisely the generalized matroid lattices which were introduced in [7].

A crucial role in our paper is played by the following theorem, the proof of which is given here in order to make the presentation more self-contained:

**THEOREM 1** (cf. [8, Lemma 1]). *Let  $L$  be a lattice of finite length. If  $c < d$  ( $c, d \in L$ ) then there exists a join-irreducible element  $u \in J(L)$  such that  $u \leq d$ ,  $u \not\leq c$  and  $u \wedge c = u'$ .*

PROOF. If  $d \in J(L)$ , then put  $u = d$ . Let now  $d \notin J(L)$ . Consider the set of all  $v \in J(L)$  which have the property

$$(1) \quad v < d \quad \text{and} \quad v \not\leq c.$$

It is obvious that this set is not empty. Choose an element  $u \in J(L)$  which is minimal with respect to property (1). From

$$u < d \quad \text{and} \quad u \not\leq c$$

it follows that

$$u \wedge c \leq u'.$$

We show that equality holds. From the assumption

$$(2) \quad u \wedge c < u'$$

we get the existence of an element  $u_* \in J(L)$  having the properties

$$(3) \quad u_* \leq u'$$

and

$$(4) \quad u_* \not\leq u \wedge c.$$

Relation (3) implies

$$(5) \quad u_* \leq u' < u < d$$

and relation (4) implies

$$(6) \quad u_* \not\leq c$$

(namely,  $u_* \leq c$  yields together with  $u_* < u$  that  $u_* \leq u \wedge c$  which contradicts (4)). Now (5) and (6) together contradict the minimality of  $u$  with respect to property (1). Hence our assumption (2) was false, i.e. we have

$$u \wedge c = u'$$

which was to be proved.

Let us remark that the preceding theorem was also implicitly used in the proof of the main result of [9].

### 3. A characterization of generalized matroid lattices

First we prove

THEOREM 2. Let  $L$  be a lattice of finite length. Then

$$(*) \quad u \in J(L) \Rightarrow (u, b)M^* \quad (\forall b \in L)$$

implies that  $L$  is (upper) semimodular.

PROOF. We have to show that  $L$  satisfies the neighbourhood condition

$$(N) \quad a \wedge b < a \Rightarrow b < a \vee b \quad (a, b \in L).$$

Assume therefore that

$$(7) \quad a \wedge b < a$$

holds. By Theorem 1 there exists a join-irreducible element  $u \in J(L)$  such that

$$u \leq a, \quad u \not\leq a \wedge b \quad \text{and} \quad u' = u \wedge (a \wedge b) = u \wedge b$$

hold. We observe that then

$$a = (a \wedge b) \vee u \quad \text{and} \quad a \vee b = (a \wedge b) \vee u \vee b = b \vee u.$$

Suppose now that there exists an element  $c \in L$  such that

$$(***) \quad b < c < a \vee b = u \vee b.$$

From this we get

$$c = c \wedge (u \vee b)$$

and  $u' = b \wedge u = c \wedge u$ , that is,

$$b = (c \wedge u) \vee b.$$

This situation is visualized in Fig. 2:

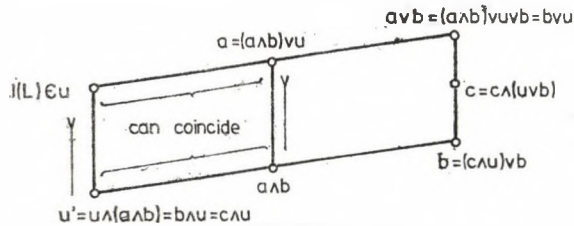


Fig. 2

Thus we get that

$$(c \wedge u) \vee b = b < c = c \wedge (u \vee b)$$

i.e. we have

$$(u, b) \overline{M^*}.$$

Therefore our assumption  $(***)$  contradicts  $(**)$ . Hence there exists no element  $c$  such that

$$b < c < a \vee b.$$

Therefore (7) yields

$$(8) \quad b < a \vee b,$$

i.e.  $L$  is (upper) semimodular which was to be proved.

The lattice in Fig. 3 shows that (upper) semimodularity does not imply

$$u \in J(L) \Rightarrow (u, b) \overline{M^*} \quad (\forall b \in L).$$

In connection with this let us remark that the lattice of Fig. 3 was used in order to characterize strongness (which was first introduced in [1]) in semimodular lattices of finite length in the following way:

THEOREM 3 (cf. [2, Theorem 6] and [6, Theorem 3]). *Let  $L$  be a semimodular lattice of finite length. Then the following two conditions are equivalent:*

- (i)  $L$  is strong;
- (ii)  $L$  has no sublattice isomorphic to the lattice in Fig. 3.

COROLLARY 4. *Let  $L$  be a lattice of finite length. Then*

$$u \in J(L) \Rightarrow (u, b)M^* \quad (\forall b \in L)$$

*implies strong semimodularity.*

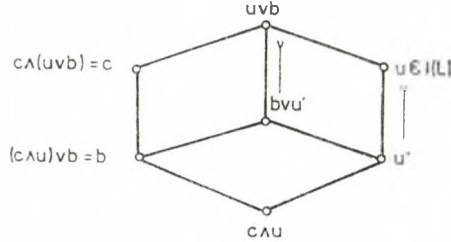


Fig. 3

PROOF.  $(u, b)M^*$  for all  $u \in J(L)$  and for all  $b \in L$  implies semimodularity by Theorem 2. If  $L$  were not strongly semimodular, then  $L$  has by Theorem 3 a sublattice isomorphic to the lattice in Fig. 3. This, however, means that  $(u, b)M^*$  does not hold and the corollary is proved.

On the other hand, even strong semimodularity does not imply  $(u, b)M^*$  for all  $u \in J(L)$  and for all  $b \in L$  which is seen from the lattice in Fig. 4:

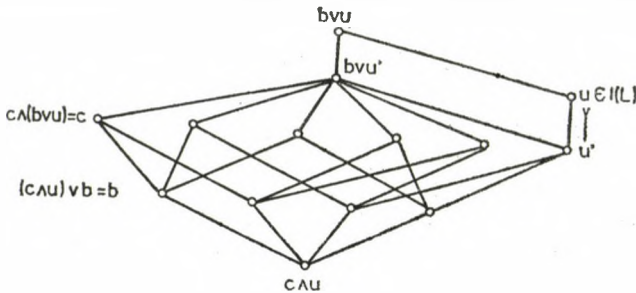


Fig. 4

It is easily checked that this lattice is strongly semimodular and that  $(u, b)M^*$  does not hold.

Using the following lemma we shall now see that the lattices of finite length in which arbitrary join-irreducible elements and arbitrary elements form dual-modular pairs are precisely the generalized matroid lattices of [7].

LEMMA 5 (cf. [5, Lemma 1.2, p. 1]). *Let  $L$  be a lattice and  $a \in L$ . Then  $(a, x)M$  ( $\forall x \in L$ ) holds if and only if  $(a, x)M^*$  ( $\forall x \in L$ ).*

We are now in a position to prove our main result:

**THEOREM 6.** *Let  $L$  be a lattice of finite length. Then the following two conditions are equivalent:*

(i)  $[u \wedge b, u] \cong [b, b \vee u]$  holds for all  $u \in J(L)$  and for all  $b \in L$ , that is,  $L$  is a generalized matroid lattice in the sense of [7];

(ii)  $(u, b)M^*$  holds for all  $u \in J(L)$  and for all  $b \in L$ , i.e. arbitrary join-irreducible elements and arbitrary elements form dual-modular pairs.

**PROOF.** (i) $\Rightarrow$ (ii): Let  $L$  be a generalized matroid lattice. From [4, Theorem 2] it follows that then  $(u, b)M^*$  holds for all  $u \in J(L)$  and for all  $b \in L$ .

(ii) $\Rightarrow$ (i): Let now  $(u, b)M^*$  hold for an arbitrary  $u \in J(L)$  and for an arbitrary  $b \in L$ .

By Theorem 2,  $L$  is a semimodular lattice.

From Lemma 5 it follows that  $(u, b)M^*(\forall u \in J(L), \forall b \in L)$  is equivalent to  $(u, b)M(\forall u \in J(L), \forall b \in L)$ . Since  $L$  is semimodular and of finite length,  $L$  is also  $M$ -symmetric, that is,  $(u, b)M$  implies  $(b, u)M$  (cf. e.g. [3, Theorem 9, p. 117]).

Now  $(b, u)M$  and  $(u, b)M^*$  yield again by [4, Theorem 2] that

$$[u \wedge b, u] \cong [b, b \vee u]$$

holds for all  $u \in J(L)$  and for all  $b \in L$ , that is,  $L$  is a generalized matroid lattice. This finishes the proof.

**COROLLARY 7.** *Let  $L$  be a finite atomistic lattice satisfying implication  $(*)$ . Then  $L$  is a classical matroid lattice.*

**COROLLARY 8.** *Each modular lattice of finite length satisfies implication  $(*)$  and is therefore a generalized matroid lattice.*

For a non-modular and non-atomistic example of a generalized matroid lattice we refer to [7].

**ACKNOWLEDGEMENTS.** The second author gratefully acknowledges support by the Natural Science and Engineering Council of Canada, Grant 214—1518. Moreover he would like to thank Professor G. Bruns (McMaster University, Hamilton, Ontario) and the Ministry of Higher Education of the German Democratic Republic for making possible a three months' stay in Canada and in the USA during spring 1982.

#### REFERENCES

- [1] FAIGLE, U., Geometries on partially ordered sets, *J. Combin. Theory. Ser. B* **28** (1980), 26—51. MR **81m**: 06020.
- [2] FAIGLE U., RICHTER, G. and STERN, M., Geometric exchange properties in lattices of finite length, *Algebra Universalis* **19** (1984), 355—365.
- [3] GRÄTZER, G., *General Lattice Theory*, Birkhäuser-Verlag, Basel und Stuttgart, 1978. MR **80c**: 06001a.
- [4] KOHL, L. and STERN, M., A characterization of certain finite lattices in which every element is a join of cycles, *Coll. Math. Soc. J. Bolyai* **33**, *Contributions to Lattice Theory*, Szeged (Hungary), 1980, 575—589.
- [5] MAEDA, F. and MAEDA, SH., *Theory of Symmetric Lattices*, Springer-Verlag, Berlin—Heidelberg—New York, 1970. MR **44**#123.

- [6] RICHTER G. and STERN, M., Strongness in (semimodular) lattices of finite length, *Wiss. Z. Martin-Luther-Univ. Halle—Wittenberg, Math.-Natur. Reihe* **33** (1984), 73—77.
- [7] STERN, M., Generalized matroid lattices, *Coll. Math. Soc. J. Bolyai* **25**, *Algebraic Methods in Graph Theory*, Szeged (Hungary), 1978, 727—748. *MR* 83a: 06007.
- [8] STERN, M., Exchange properties in lattices of finite length, *Wiss. Z. Martin-Luther-Univ. Halle—Wittenberg Math. Natur.* **31** (1982), 15—26.
- [9] STERN, M., Semimodularity in lattices of finite length, *Discrete Math.* **41** (1982), 287—293. *MR* 83m: 06010.

(Received February 14, 1984)

MARTIN-LUTHER-UNIVERSITÄT  
SEKTION MATHEMATIK  
DDR-4010 HALLE  
UNIVERSITÄTSPLATZ 6  
GERMAN DEMOCRATIC REPUBLIC





## ON THE DEFORMATIONS OF $L_1$

A. FIALOWSKI

Let  $W_1$  be the Lie algebra of the vector fields on the line with polynomial coefficients, and  $L_1$  its subalgebra, consisting of the fields, which turn into zero with their first derivatives at the origin. In this work we study the deformations of  $L_1$ . Calculation of various invariants of infinite dimensional Lie algebras is rather complicated, even in the case of affine Lie algebras (see [4]).  $L_1$  seems in some sense the next one by difficulty. The main results of the present paper, without proofs, have been published in [3].

The study of deformations of nilpotent subalgebras in the Lie algebra of vector fields was stimulated specifically by [2], where it was proved that the cohomology with trivial coefficients of  $L_1$  and similar nilpotent Lie (super) algebras of geometrical origin are invariant under certain deformations. Moreover,  $L_1$  seems to be the simplest one among the Lie algebras which have higher order obstructions to continuation of deformations, as follows from the properties of its cohomology with coefficients in the coadjoint representation, which are studied in this work.

All these show that the study of deformations of  $L_1$  has to be one of the first steps on the way to study deformations of infinite dimensional Lie algebras.

1. At first we give some facts about the cohomology of  $L_1$  with coefficients in the coadjoint representation.

In  $W_1$  choose the basis  $e_{-1}, e_0, e_1, e_2, \dots$ , where  $e_i = x^{i+1} \partial / \partial x$ . The bracket in this basis is of the form  $[e_i, e_j] = (j-i)e_{i+j}$ . Then  $L_1$  has the basis  $e_1, e_2, \dots$ . We also need the subalgebras  $L_i \subset W_1$ ; the basis in  $L_i$  ( $i \geq 0$ ) consists of the fields  $e_i, e_{i+1}, \dots$ .  $W_1$  and  $L_i$  are naturally graded, the grading (weight) of  $e_i$  equals  $i$ . This grading is inherited by the cohomology spaces with coefficients in the graded modules. Specifically,  $H^k(L_1; L_1) = \bigoplus_m H^k_{(m)}(L_1; L_1)$ .

PROPOSITION 1.  $H^k(L_1; L_1)$ ,  $k > 0$  has dimension  $2k-1$  and is generated by elements of weight  $-\frac{3k^2-k}{2} + i$ ,  $i = 1, 2, \dots, 2k-1$ . (In other words,  $H^k_{(m)}(L_1; L_1) \cong H^{k-1}_{(m)}(L_2; 1)$ ).

Specifically,  $H^1(L_1; L_1)$  is one-dimensional with weight 0;  $H^2(L_1; L_1)$  is three-dimensional and is generated by elements of weights  $-2, -3, -4$ ;  $H^3(L_1; L_1)$  is five-dimensional and is generated by elements of weights  $-7, -8, -9, -10, -11$ .

1980 *Mathematics Subject Classification*. Primary 17B56; Secondary 17B65.

*Key words and phrases*. Deformation, cohomology class, obstruction, Lie operation of Massey.

PROOF. This proposition follows from the results of [1]. Namely, the homology groups of  $L_1$  with coefficients in the modules  $\mathcal{F}_{\lambda,\mu}$ ,  $F_{\lambda,\mu}$ ,  $F'_{\lambda,\mu}$  are calculated in [1]. Here  $\mathcal{F}_{\lambda,\mu}$ ,  $\lambda, \mu \in \mathbb{C}$  is the module with the basis  $f_i$ ,  $i \in \mathbb{Z}$ ,  $e_i f_j = (j + \mu - \lambda(i+1))f_{i+j}$ ;  $F_{\lambda,\mu}$  is the submodule of  $\mathcal{F}_{\lambda,\mu}$  with the basis  $f_0, f_1, f_2, \dots$ ;  $F'_{\lambda,\mu}$  is the module, conjugate to  $F_{\lambda,\mu}$ . The adjoint representation in this notation is  $F_{1,1}$ . Remark that  $H^i(L_1, F_{1,1})$  is dual to  $H_i(L_1, F'_{1,1})'$ . By using Theorem 4.3 in [1] we get  $H_i(L_1, F'_{1,1}) = H_{i-1}(L_1, F_{-2,-1}) \oplus H_i(L_1, \mathcal{F}_{-2,-1})$ . From Theorem 4.1 [1] it follows that  $H_i(L_1, \mathcal{F}_{-2,-1}) = 0$  for each  $i$ . Using Theorem 4.2 [1] and the subsequent we receive that  $\dim H_i(L_1, F_{-2,-1}) = \dim H_i(L_2)$ . Q.e.d.

The spaces  $H^*(L_k; 1)$  are calculated in [5]. By comparing the results of those calculations with the above one we obtain the required Proposition. Cocycle  $\varphi$  representing  $H^1(L_1; L_1)$  has the form  $\varphi(e_i) = ie_i$ . The elements of  $H^1(L_1; L_1)$ , as we know, are exterior derivations. Each of these derivations define a Lie algebra, containing  $L_1$  as its ideal of codimension 1. In the present case we obtain  $L_0$ .

Let us denote the three homogeneous elements of weights  $-2$ ,  $-3$  and  $-4$  in  $H^2(L_1; L_1)$  by  $\alpha, \beta$  and  $\gamma$ , respectively. We give explicitly the cocycles  $\bar{\alpha}, \bar{\beta}, \bar{\gamma}$ , which represent the cohomology classes  $\alpha, \beta, \gamma$ :

$$\begin{aligned} \bar{\alpha}(e_2, e_3) &= 4e_3, \\ \bar{\alpha}(e_2, e_j) &= je_j \quad (j \geq 4), \quad \bar{\alpha}(e_3, e_j) = -(j-1)e_{j+1} \quad (j \geq 4), \\ \bar{\alpha}(e_i, e_j) &= 0 \quad \text{for other } i, j; \\ \bar{\beta}(e_2, e_3) &= 8e_2, \quad \bar{\beta}(e_2, e_4) = 4e_3, \quad \bar{\beta}(e_3, e_4) = -10e_4, \\ \bar{\beta}(e_2, e_j) &= (j+1)e_{j-1}, \quad \bar{\beta}(e_3, e_j) = -2je_j, \quad \bar{\beta}(e_4, e_j) = (j-1)e_{j+1} \quad \text{for } j \geq 5, \\ (1) \quad \bar{\beta}(e_i, e_j) &= 0 \quad \text{for other } i, j; \\ \bar{\gamma}(e_2, e_3) &= 14e_1, \quad \bar{\gamma}(e_2, e_5) = 8e_3, \quad \bar{\gamma}(e_3, e_4) = -24e_3, \\ \bar{\gamma}(e_3, e_5) &= -16e_4, \quad \bar{\gamma}(e_4, e_5) = 18e_5, \\ \bar{\gamma}(e_2, e_j) &= (j+2)e_{j-2}, \quad \bar{\gamma}(e_3, e_j) = -3(j+1)e_{j-1} \quad \left. \vphantom{\bar{\gamma}(e_2, e_j)} \right\} \quad \text{for } j \geq 6, \\ \bar{\gamma}(e_4, e_j) &= 3je_j, \quad \bar{\gamma}(e_5, e_j) = -(j-1)e_{j+1} \\ \bar{\gamma}(e_i, e_j) &= 0 \quad \text{for other } i, j. \end{aligned}$$

Indeed, it is easy to prove that every two-dimensional cohomology class of  $L_1$  with coefficients in  $L_1$  is represented by a single cocycle  $\omega$ , which turns into 0 at the field  $e_1$  (i.e.  $\omega(e_1, e_j) = 0$  for each  $j$ ). Then for  $\omega(e_i, e_j)$  with  $1 < i < j$  we can write a system of equations, which has a unique solution (up to a constant multiplier). By giving it explicitly we give another proof of the part of Proposition 1 related to the two-dimensional cohomology.

REMARK. A simpler formula for the cohomology class  $\alpha$  may be obtained in the following way. Let  $e \in H^1(L_1; 1)$  be an element of weight  $-2$  and  $v \in H^1(L_1; L_1)$  the cohomology class of the cocycle  $e_i \mapsto ie_i$ ; then  $\alpha = ev$ . Classes  $\beta$  and  $\gamma$  have no simple description of this kind.

## 2. Obstructions of deformations of $L_1$ .

To find the versal family of  $L_1$  we use the methods of [6]. In this work the Lie operations of Massey are defined and the next process of finding the tangent cone to the base of the versal deformation is proposed. Let  $K_1 \subset H^2(L_1; L_1)$  (where  $H^2(L_1; L_1)$  is the tangent space to the base of the versal deformation) be the cone, consisting of all  $\xi$  such that  $\langle \xi, \xi \rangle = 0$ . Let further  $K_2$  be the cone, consisting of those elements  $\xi \in K_1$  that  $\langle \xi, \xi, \xi \rangle \geq 0$ , and in general,  $K_i = \{\xi \in K_{i-1} | \langle \xi, \dots, \xi \rangle_{i+1} \geq 0\}$ . The tangent cone to the base of the versal family is the intersection of all  $K_i$ .

In general, let  $\mathcal{L}$  be a Lie algebra,  $[\ , \ ]_t$  the deformation of the bracket. Let us expand this deformation into a Taylor series:

$$[x, y]_t = [x, y]_0 + t\omega_1(x, y) + t^2\omega_2(x, y) + \dots, \quad x, y \in \mathcal{L}.$$

The Jacobi identity is:

$$\begin{aligned} [[x, y]_t, z]_t &= [[x, y]_0 + t\omega_1(x, y) + t^2\omega_2(x, y) + \dots, z]_t = [[x, y]_0, z]_0 + \\ &+ t\omega_1([x, y]_0, z) + t^2\omega_2([x, y]_0, z) + t[\omega_1(x, y), z]_0 + t^2\omega_1(\omega_1(x, y), z) + \\ &+ t^2[\omega_2(x, y), z]_0 + \dots = [[x, z]_t, y]_t + [x, [y, z]_t]_t. \end{aligned}$$

From this equation it follows that  $\omega_1$  is a 2-cocycle of  $\mathcal{L}$  with coefficients in the coadjoint representation and also that the 3-cochain

$$\omega_1(\omega_1(x, y), z) - \omega_1(\omega_1(x, z), y) + \omega_1(\omega_1(y, z), x)$$

is coboundary of  $\omega_2$ . Let  $\alpha \in H^2(\mathcal{L}; \mathcal{L})$  be the cohomology class of  $\omega_1$ . We actually prove that the Lie square  $\langle \alpha, \alpha \rangle$  equals 0. For the continuation of the infinitesimal deformation  $[x, y] + t\omega_1[x, y]$  to  $\text{spec } k[t]/t^3$  it is necessary that  $\langle \alpha, \alpha \rangle = 0$ .

If  $\mathcal{L} \cong L_1$  then, as easy to see,  $\langle \alpha, \alpha \rangle = \langle \beta, \beta \rangle = 0$ , because the weights of  $\langle \alpha, \alpha \rangle$  and  $\langle \beta, \beta \rangle$  are  $-4$  and  $-6$ , and there is no such three-dimensional cohomology class. Similarly, by consideration of the dimensions we have  $\langle \alpha, \beta \rangle = \langle \alpha, \gamma \rangle = 0$ . The calculation of the remaining relations is less evident.

**PROPOSITION 2.** *The Lie products  $\langle \gamma, \gamma \rangle$ ,  $\langle \beta, \gamma \rangle$  and the Massey Lie cube  $\langle \beta, \beta, \beta \rangle$  are nontrivial, while  $\langle \alpha, \alpha, \dots, \alpha \rangle \geq 0$  for all  $i$ .*

**PROOF.** To prove  $\langle \gamma, \gamma \rangle \neq 0$  substitute  $\bar{\gamma}$  into (1). The three-dimensional cocycle we obtain is not cohomologous to 0, as its value on the cohomology class of weight  $-8$  is different from zero. Direct computation shows that  $\langle \bar{\gamma}, \bar{\gamma} \rangle$  has nonzero value on the class of weight  $-8$ .

The fact that  $\langle \beta, \gamma \rangle \neq 0$  we may prove in the same way; however, we give another proof for this.

Now  $\langle \alpha, \alpha, \dots, \alpha \rangle \geq 0$  follows from the fact that there exists a deformation with infinitesimal deformation  $\alpha$ .

The most laborious part of the proof is to show that  $\langle \beta, \beta, \beta \rangle \neq 0$ . This is equivalent to the fact that for any Lie algebra over  $k[t]/t^4$  with the basis  $e_i$  the bracket cannot be of the following form:

$$[e_i, e_j]_t = (j-i)e_{i+j} + t\beta(e_i, e_j)e_{i+j-3} + t^2\kappa_1(e_i, e_j)e_{i+j-6} + t^3\kappa_2(e_i, e_j)e_{i+j-9}.$$

Here  $\kappa_1(e_i, e_j)$  and  $\kappa_2(e_i, e_j)$  may be defined step by step ( $i+j=1, 2, \dots$ ) from the system of equations following from the Jacobi identity. For  $i+j=12$  we get a contradiction. Jacobi identities of weight 11 define an algebra  $L$  of dimension 11.  $H^2(L; 1)$  contains no elements of weight 12. This can be explained in the following way.

The Lie algebra  $\tilde{L}$  with the basis  $e_1, \dots, e_{11}$  and the bracket  $[e_i, e_j] = (j-i)e_{i+j}$  has, similarly to  $L_1$ , a three-dimensional cohomology class of weight 12 and a two-dimensional cohomology class of weight 12. But  $L$  does not have any of these. It has a filtration, and the corresponding coadjoint graded algebra is  $\tilde{L}$ . Hence there exists a spectral sequence with the first term  $H^*(\tilde{L}; 1)$  converging to  $H^*(L; 1)$ . In this spectral sequence the differential of the two-dimensional class of weight 12 is a three-dimensional class of weight 12.

It remains to verify that  $\langle \beta, \gamma \rangle \neq 0$ . Here equality would imply that  $L_1$  has a deformation over  $\text{spec } k[t_1, t_2]/(t_1^2, t_2^2)$  such that in the deformed algebra

$$[e_i, e_j]_t = (j-i)e_{i+j} + t_1 \bar{\beta}(e_i, e_j)e_{i+j-3} + t_2 \bar{\gamma}(e_i, e_j)e_{i+j-4} + t_1 t_2 \kappa(e_i, e_j)e_{i+j-7}.$$

Straight calculation shows that such an algebra cannot exist. Namely, the numbers  $\kappa_{i,j} = \kappa(e_i, e_j)$  can be defined step by step. At the 12-th step we get a system of equations for  $\kappa_{ij}$ , which has no solution.

Let us define now three deformations of the Lie algebra structure in  $L_1; [\ , \ ]_t^1, [\ , \ ]_t^2, [\ , \ ]_t^3$ .

$$\begin{aligned} [e_i, e_j]_t^1 &= (j-i)(e_{i+j} + te_{i+j-1}); \\ (2) \quad [e_i, e_j]_t^2 &= \begin{cases} (j-i)e_{i+j}, & \text{if } i, j > 1, \\ (j-1)e_{j+1} + te_j, & \text{if } i = 1, \ j > 1; \end{cases} \\ [e_i, e_j]_t^3 &= \begin{cases} (j-i)e_{i+j}, & \text{if } i, j \neq 2, \\ (j-2)e_{j+2} + te_j, & \text{if } i = 2, \ j \neq 2. \end{cases} \end{aligned}$$

These three Lie algebra families may be realized inside the Lie algebra  $L_0$ . By the first deformation  $e_i$  deforms into  $e_i + te_{i-1}$ ,  $i > 0$ . In other words,  $L_1$  deforms into the Lie algebra of vector fields on the line, vanishing at the origin and  $t$ . By the second deformation the  $e_i$ 's,  $i > 1$  are unchanged and  $e_1$  deforms into  $e_1 + te_0$ . Analogously, by the third deformation  $e_2$  turns into  $e_2 + te_0$  and the remaining elements are unchanged. We show now that no other deformations exist.

From Proposition 2 it follows that the tangent cone to the base of versal deformation of  $L_1$  is a curve of degree 3, generated by  $\alpha$ . The versal deformation consists of three curves, tangential to  $\alpha$ . These three curves correspond to the three Lie algebra families.

It follows then

**THEOREM.** *Each deformation of the Lie algebra structure in  $L_1$  may be obtained by substitution of the variable from one of the three families  $[\ , \ ]_t^1, [\ , \ ]_t^2, [\ , \ ]_t^3$  (see formulas (2)).*

PROOF. It is sufficient to verify that the three families are nontrivial and pairwise nonisomorphic. Indeed, the Lie algebras with the brackets  $[\cdot, \cdot]_t, [\cdot, \cdot]_{t'}, [\cdot, \cdot]_{t''}$ ,  $i \neq j, t', t'' \neq 0$  are not isomorphic neither with each other nor with  $L_1$ . (Notice that the algebras with  $[\cdot, \cdot]_t, [\cdot, \cdot]_{t'}, [\cdot, \cdot]_{t''}$  are isomorphic for  $t', t'' \neq 0$ .) Let  $L_1(i)$  be an algebra from the  $i$ -th family. First,  $L_1(i) \not\cong L_1$  as  $L_1$  is a nilpotent algebra, while  $L_1(i)$ , for all  $i$ , is only solvable. Further, the maximal nilpotent subalgebra in  $L_1(1)$  has codimension two, and in  $L_1(2)$  and  $L_1(3)$  codimension one. Finally, the algebra  $[L_1(3), L_1(3)]$  has two generators, while  $[L_1(2), L_1(2)]$  has three. So  $L_1(1)$ ,  $L_1(2)$  and  $L_1(3)$  are neither isomorphic to each other, nor to  $L_1$ .

REMARKS. 1. From the Theorem it follows that all the nontrivial deformations of  $L_1$  lose the grading. In the class of graded Lie algebras  $L_1$  has no nontrivial deformations.

2. Let us study the change of grading by deformations of  $L_1$  in detail. Let  $L$  be a Lie algebra, lying in one of the three families. Choose in  $L$  the following basis:  $\{\bar{e}_1, \bar{e}_2, \dots\}$ ,  $\bar{e}_1 = e_1$ ,  $\bar{e}_2 = e_2 + \alpha e_1$ ;  $\bar{e}_3, \bar{e}_4, \dots$  are defined so that  $[\bar{e}_1, \bar{e}_i] = (i-1)\bar{e}_{i+1}$  holds. In the new basis we have:

$$[\bar{e}_i, \bar{e}_j] = (j-i)\bar{e}_{i+j} + \omega_1(i, j)\bar{e}_{i+j-1} + \omega_2(i, j)\bar{e}_{i+j-2} + \dots$$

Here  $\omega_1$  is a 2-cocycle of  $L_1$  with coefficients in  $L_1$  of weight  $-1$  such that  $\omega_1(e_1, e_j) = 0$ . Such a cocycle is represented as a differential  $\omega = dv$ . It is easy to see that  $v(e_i) = 0, i \neq 2$  and  $v(e_2)$  is proportional to  $e_1$ . Hence the parameter  $\alpha$  may be chosen so that  $\omega_1 = 0$ . Let us call a basis in  $L$  canonical if  $\omega_1 = 0$ .

Notice that  $\omega_2, \omega_3$  are cocycles of  $L_1$  and  $d\omega_4 = [\omega_1, \omega_1]$ . Cocycles  $\omega_2, \omega_3$  are proportional to those in (1). Choose in the three families three algebras —  $L^1, L^2, L^3$  — such that  $\omega_2 = \omega_2(L^1)$  exactly coincides with the first cocycle in (1). For the first and third family we have  $\omega_3 = 0$ , but for the second one  $\omega_3 \neq 0$ . Further

$$d(\omega_4(L^2) - \omega_4(L^1)) = 0 \quad \text{and} \quad d(\omega_4(L^3) - \omega_4(L^1)) = 0.$$

Cocycles  $\bar{u}_1 = \omega_4(L^2) - \omega_4(L^1)$  and  $\bar{u}_2 = \omega_4(L^3) - \omega_4(L^1)$  represent nontrivial cohomology classes of weight  $-4$  from  $H^2(L_1; L_1)$ . These cocycles are proportional to the third one in (1). Consequently, the quotient of these cocycles with coefficients  $\mu_1$  and  $\mu_2$ , respectively,  $\mu = \mu_1/\mu_2$  is a number, which may be easily computed. This number occurs as one of the numerical characteristics of the versal deformation of  $L_1$ . It would be interesting to connect this number with other characteristics of the deformation.

The author thanks B. L. Feigin and D. B. Fuchs for their helpful comments.

## REFERENCES

- [1] FEIGIN, B. L. and FUCHS, D. B., Homology of the Lie algebra of the vector fields on the line, *Funkcional. Anal. i Priložen.* **14** (1980), 45—60. *MR 82b*: 17017.
- [2] FEIGIN, B. L. and RETAH, V. S., On the cohomology of certain Lie algebras and superalgebras of vector fields, *Uspehy Mat. Nauk* **37** (1982), 233—234. *MR 83m*: 58083.
- [3] FIALOWSKI, A., Deformations of the Lie algebra of vector fields on the line, *Uspehy Mat. Nauk* **38** (1983), 201—202 (Russian). *MR 84m*: 17011.
- [4] FIALOWSKI, A., Deformations of nilpotent Kac—Moody algebras, *Studia Sci. Math. Hungar.* **19** (1984),
- [5] GONCHAROVA, L. V., Cohomology of the Lie algebra of formal vector fields on the line, *Funkcional. Anal. i Priložen.* **7** (1973), 6—14. *MR 49* #4058a.
- [6] RETAH, V. S., Massey operations in Lie superalgebras and deformations of complex analytic manifolds, *Funkcional. Anal. i Priložen.* **11** (1977), 88—89. *MR 58* #22670.

(Received March 7, 1984)

VILLÁNYI ÚT 103  
H-1118 BUDAPEST  
HUNGARY



# AN ARCSINE-LAW FOR THE OSCILLATING RANDOM WALK

PETER REIMNITZ<sup>1</sup>

## Summary

We are given a Markov chain  $(Z_n)_{n=0}^{\infty}$  defined by

$$Z_0 = z, \quad z \in \mathbb{R}$$

$$Z_{n+1} = Z_n + \begin{cases} Y_{1,n+1} & \text{if } Z_n \geq 0, \\ Y_{2,n+1} & \text{if } Z_n < 0. \end{cases}$$

Here,  $(Y_{1,i})_{i=1}^{\infty}$  and  $(Y_{2,i})_{i=1}^{\infty}$  are two independent sequences of independent, identically distributed random variables with

$$E[Y_{1,1}] = E[Y_{2,1}] = 0 \quad \text{and} \quad E[Y_{1,1}^2] = \sigma_1^2 < \infty, \quad E[Y_{2,1}^2] = \sigma_2^2 < \infty.$$

It is shown that

$$\frac{1}{n} \sum_{i=1}^n 1_{(Z_i \geq 0)} \quad \left( 1_A(\omega) = \begin{cases} 1 & \text{if } \omega \in A \\ 0 & \text{if } \omega \notin A \end{cases} \right)$$

converges in distribution towards a random variable with density (with respect to the Lebesgue measure):

$$g(y) = \frac{R}{\pi} (y(1-y))^{-1/2} (y + (1-y)R^2)^{-1} 1_{(0,1)}(y).$$

## 1. Introduction

We are given a Markov chain  $(Z_n)_{n=0}^{\infty}$  defined by

$$Z_0 = z, \quad z \in \mathbb{R};$$

$$(1) \quad Z_{n+1} = Z_n + \begin{cases} Y_{1,n+1} & \text{if } Z_n \geq 0, \\ Y_{2,n+1} & \text{if } Z_n < 0. \end{cases}$$

Here,  $(Y_{1,i})_{i=1}^{\infty}$  and  $(Y_{2,i})_{i=1}^{\infty}$  are two independent sequences of independent, identically distributed random variables with  $E[Y_{1,1}] = E[Y_{2,1}] = 0$  and  $E[Y_{1,1}^2] = \sigma_1^2 < \infty$ ,  $E[Y_{2,1}^2] = \sigma_2^2 < \infty$ .

Define a sequence of random variables  $(N_n)_{n=0}^{\infty}$  through

$$(2) \quad N_n = \sum_{i=0}^n 1_{(Z_i \geq 0)},$$

<sup>1</sup> Research was supported by "Deutsche Forschungsgemeinschaft" under SFB 72.  
1980 *Mathematics Subject Classification*. Primary 60G50; Secondary 60F05.  
*Key words and phrases*. Arcsine-law, oscillating random walk.

where

$$1_A(\omega) = \begin{cases} 1 & \text{if } \omega \in A \\ 0 & \text{if } \omega \notin A \end{cases}$$

( $A$  an arbitrary measurable set).

We will show that the distribution of  $N_n/n$  converges against a limiting distribution function with density (with respect to the Lebesgue measure):

$$(3) \quad g(y) = \frac{R}{\pi} (y(1-y))^{-1/2} (y + (1-y)R^2)^{-1} 1_{(0,1)}(y),$$

where  $R = \sigma_2/\sigma_1$ .

This result generalizes previous results for lattice distributions, where it was additionally assumed that  $P(Y_{1,1} = -1 | Y_{1,1} < 0) = 1$  and  $P(Y_{2,1} = +1 | Y_{2,1} > 0) = 1$  (skip-free distributions) (see Lamperti (1958) and Reimnitz (1977)). Keilson and Wellner (1977) derived an arcsine-law for the oscillating Brownian motion.

## 2. Auxiliary results

We first derive the generating function  $\Phi(t, \varrho) = \sum_{m=0}^{\infty} E[\varrho^{N_m}] t^m$ . For this, we observe that the functions  $\varphi_n(s, \varrho) = E[e^{sZ_n} \varrho^{N_n}]$  obey the following recursion formula:

$$(4) \quad \varphi_{n+1}(s, \varrho) = \varrho^{\hat{v}}(s) \varphi_n^+(s, \varrho) + \hat{\mu}(s) \varphi_n^-(s, \varrho)$$

where

$$\hat{v}(s) = E[\exp(sY_{1,1})], \quad \hat{\mu}(s) = E[\exp(sY_{2,1})],$$

$$\varphi_n^+(s, \varrho) = E[\exp(sZ_n) \varrho^{N_n} 1_{(Z_n \geq 0)}], \quad \varphi_n^-(s, \varrho) = \varphi_n(s, \varrho) - \varphi_n^+(s, \varrho), \quad s \in \mathbb{C}$$

(the formula holds for those  $s \in \mathbb{C}$  for which  $\hat{v}(s)$  and  $\hat{\mu}(s)$  exist).

Formula (4) is completely analogous to formula (3.8) found in Kemperman (1974). Therefore, we can use corresponding results (formulae (3.14) and (3.16)).

$$(5) \quad \begin{aligned} Q^-(t, \varrho, s) &= (\hat{\delta}_z \lambda)^- \exp(\mathcal{L}_\mu^- - \mathcal{L}_{\varrho v}^-) \\ Q^+(t, \varrho, s) &= (\hat{\delta}_z \lambda)^+ \exp(\mathcal{L}_{\varrho v}^+ - \mathcal{L}_\mu^+) + Q_0 \exp(\mathcal{L}_{\varrho v}^+ - \mathcal{L}_\mu^+ - 1), \end{aligned}$$

here  $Q^\pm(t, \varrho, s)$  are the functions:

$$Q^+(t, \varrho, s) = \sum_{m=0}^{\infty} E[\varrho^{N_m} \exp(sZ_m) 1_{(Z_m \geq 0)}] t^m,$$

$$Q^-(t, \varrho, s) = \sum_{m=0}^{\infty} E[\varrho^{N_m} \exp(sZ_m) 1_{(Z_m < 0)}] t^m,$$

$Q_0$  is zero for  $z \neq 0$  and one for  $z = 0$ ;

the functions  $\mathcal{L}_\pi^\pm(t, s)$  are defined to be

$$\mathcal{L}_\pi^+(t, s) = \sum_{n=1}^{\infty} \frac{t^n}{n} (\hat{\pi}^n)^+(s) \quad ((\hat{\pi}^n)^+(s) = \int_{0^-}^{\infty} e^{sx} d\pi^{*n}(x)),$$

$$\mathcal{L}_\pi^-(t, s) = \sum_{n=1}^{\infty} \frac{t^n}{n} (\hat{\pi}^n)^-(s) \quad ((\hat{\pi}^n)^-(s) = \hat{\pi}^n(s) - (\hat{\pi}^n(s))^+),$$

$\pi$  is an arbitrary probability and  $\pi^{*n}$  denotes the  $n$ -fold convolution of  $\pi$ ;  $\lambda$  is defined to be

$$\lambda(t, s, \varrho) = \exp(\mathcal{L}_\mu^+(t, s) + \mathcal{L}_{\varrho^v}^-(t, s));$$

finally

$$(\hat{\delta}_z \lambda)^+ = \int_{0^-}^{\infty} e^{sx} d\delta_z * \lambda(x), \quad (\hat{\delta}_z \lambda)^- = e^{sz} \lambda(s) - (\hat{\delta}_z \lambda)^+$$

here  $\delta_z * \lambda(x)$  is the distribution function of the (uniquely determined) measure with  $\int e^{sx} d\delta_z * \lambda(x) = e^{sz} \lambda(s)$ . Of course,  $\Phi(\varrho, t) = Q^-(t, \varrho, 0) + Q^+(t, \varrho, 0)$ .

In the next section we will derive explicit formulae for  $Q^-(t, \varrho, s)$  and  $Q^+(t, \varrho, s)$  for special cases of distribution functions of  $Y_{1,1}$  and  $Y_{2,1}$ .

### 3. Explicit formulae for the generating functions in special cases

In this section we assume the Laplace transforms of  $Y_{1,1}$  and  $Y_{2,1}$  to be of the forms:

$$\begin{aligned} E[\exp(sY_{1,1})] &= \hat{\nu}(s) \\ &= \sum_{i=1}^n \frac{n_i}{s+b_i} + E[\exp(sY_{1,1}) 1_{(Y_{1,1} \geq 0)}] \end{aligned}$$

(6)

$$\begin{aligned} E[\exp(sY_{2,1})] &= \hat{\mu}(s) \\ &= \sum_{i=1}^m \frac{m_i}{s-a_i} + E[\exp(sY_{2,1}) 1_{(Y_{2,1} < 0)}], \end{aligned}$$

hereby  $n, m \in \mathbb{N}$ ,  $n_i, m_i \in \mathbb{R}$ ,  $b_i, a_i \in \mathbb{R}^+$ . Furthermore, we assume  $\hat{\nu}(s)$  and  $\hat{\mu}(s)$  to exist in an open strip around the imaginary axis.

Using the Wiener—Hopf decomposition method, we obtain similarly to the results in Kemperman (1961):

$$\begin{aligned} \exp(-\mathcal{L}_{\varrho^v}^-(t, s)) &= \prod_{i=1}^n ((s+v_i(\varrho t))/(s+b_i)), \\ \exp(-\mathcal{L}_{\varrho^v}^+(t, s)) &= (1-t\varrho\hat{\nu}(s)) \prod_{i=1}^n ((s+b_i)/(s+v_i(\varrho t))), \\ \exp(-\mathcal{L}_\mu^+(t, s)) &= \prod_{i=1}^m ((s-u_i(t))/(s-a_i)), \\ \exp(-\mathcal{L}_\mu^-(t, s)) &= (1-t\hat{\mu}(s)) \prod_{i=1}^m ((s-a_i)/(s-u_i(t))), \end{aligned}$$

(7)

here,  $-v_i(t)$  are the  $n$  (counting multiplicities) roots of the equation  $1-t\hat{\nu}(s)=0$  in the negative halfplane ( $|t|<1$ ), similarly  $u_i(t)$  are the  $m$  (counting multiplicities) roots of the equation  $1-t\hat{\mu}(s)=0$  in the positive halfplane ( $|t|<1$ ). The roots  $-v_i(t)$  and  $u_i(t)$  ( $i=1, 2, \dots, n$ ) are assumed to be ordered such that for all  $i < j$ :  $u_i(t) \leq u_j(t)$  and  $v_i(t) \leq v_j(t)$  ( $|t|<1$ ).

We then obtain

$$(8) \quad \lambda(t, s, \varrho) = \sum_{i=1}^m ((s-a_i)/(s-u_i(t))) \prod_{j=1}^n ((s+b_j)/(s+v_j(\varrho t))).$$

It is easy to verify

$$(9) \quad (\delta_z \lambda)^- = \sum_{j=1}^n \exp(-v_j(\varrho t)z) A_j(t, \varrho) (s+v_j(\varrho t))^{-1},$$

$$(\delta_z \lambda)^+ = e^{sz} \lambda - (\delta_z \lambda)^-,$$

where

$$A_j(t, \varrho) = (-v_j(\varrho t) + b_j) \prod_{i \neq j}^n ((-v_j(\varrho t) + b_i)/(-v_j(\varrho t) + v_i(\varrho t))) \times \\ \times \prod_{i=1}^m ((v_j(\varrho t) + a_i)/(v_j(\varrho t) + u_i(t))).$$

Equations (5), (7) and (9) already give an explicit formula for  $\Phi(t, \varrho)$

$$(10) \quad \Phi(t, \varrho) = e^{-1} Q_0 (1-t\varrho)^{-1} \prod_{i=1}^n (v_i(\varrho t)/b_i) \prod_{i=1}^m (u_i(t)/a_i) + \\ + (1-t)^{-1} \prod_{i=1}^m (u_i(t)/a_i) \prod_{i=1}^n (v_i(\varrho t)/b_i) \cdot \\ \cdot \sum_{j=1}^n \exp(-v_j(\varrho t)z) v_j(\varrho t)^{-1} A_j(t, \varrho) + (1-\varrho t)^{-1} + \\ + (1-\varrho t)^{-1} \prod_{i=1}^m (u_i(t)/a_i) \prod_{i=1}^n (v_i(t)/b_i) \cdot \\ \cdot \sum_{j=1}^n \exp(-v_j(\varrho t)z) v_j(\varrho t)^{-1} A_j(t, \varrho).$$

From Corollary 1 of Appendix A we find:

$$\lim_{m \rightarrow \infty} E[(1-u(N_m/m))^{-1}] = \lim_{t \rightarrow 1} (1-t) \Phi(t, \varrho^{1-t}),$$

where  $u = \log \varrho$ .

From this result we immediately obtain our main theorem for the special class of (underlying) distribution functions fulfilling (6):

**THEOREM 1.** *Let the distribution functions of  $Y_{1,1}$  and  $Y_{2,1}$  fulfill condition (6) and  $E[Y_{1,1}] = E[Y_{2,1}] = 0$ ,  $E[Y_{1,1}^2] = \sigma_1^2$ ,  $E[Y_{2,1}^2] = \sigma_2^2$ ; further let  $N_n$  be defined by (2), then*

$$\lim_{n \rightarrow \infty} \mathcal{L}(N_n/n) = \mathcal{L}(Y),$$

where  $Y$  is a random variable with density  $g(y)$  (with respect to the Lebesgue measure) and

$$g(y) = \pi^{-1} R(y(1-y))^{-1/2} (y + (1-y)R^2)^{-1} 1_{[0,1]}(y), \quad R = \sigma_2/\sigma_1$$

( $\mathcal{L}(X)$  denotes the distribution function of the random variable  $X$ ).

PROOF.  $E[1/(v - (N_m/m))]$  is the Stieltjes-transform of the random variable  $N_m/m$ . Similarly to the situation by Laplace transforms, the convergence of the Stieltjes-transforms of a sequence of random variables to a function  $g(v)$  is equivalent to convergence of the distributions of this sequence of random variables to a (possibly defective) distribution function with Stieltjes transform  $g(v)$ , see Widder (1972).

The density of a random variable is the imaginary part of the Stieltjes-transform (Widder (1972)). Therefore we are finished if we can show:

$$(11) \quad \lim_{t \rightarrow 1} (1-t) \Phi(t, e^{u(1-t)}) = (1 - R\sqrt{1-u})^{-1} + R^2(R^2(1-u) + R\sqrt{1-u})^{-1}.$$

To prove (11) we first observe:

$$\lim_{t \rightarrow 1} \frac{1-t}{1 - e^{(1-t)t}} = \frac{1}{1 - \log e},$$

$$\lim_{t \rightarrow 1} \frac{u_1(t)}{\sqrt{1-t}} = \sqrt{\frac{2}{\sigma_2^2}}$$

and

$$\lim_{t \rightarrow 1} \frac{v_1(t)}{\sqrt{1-t}} = \sqrt{\frac{2}{\sigma_1^2}}.$$

The result then follows by simple algebra from (10). ■

#### 4. Extensions of the main result to arbitrary non lattice distributions with finite second moment

The family of distributions fulfilling (6) is dense in the space of all distributions over  $\mathbf{R}$  in the usual topology of weak convergence, see Appendix B. Therefore, we will approximate arbitrary distribution functions by distribution functions fulfilling (6).

To indicate dependency on underlying distributions we will henceforward write

$$\Phi_{\mu, v}(t, \varrho) = \Phi(t, \varrho) \quad \text{and} \quad \lambda_{\mu, v}(t, \varrho, s) = \lambda(t, \varrho, s).$$

Essentially we have to prove

$$\lim_{n \rightarrow \infty} (1-t) \Phi_{\mu_n, v_n}(t, \varrho) = (1-t) \Phi_{\mu, v}(t, \varrho) \quad \text{for all } |t| < 1,$$

and

$$\lim_{n \rightarrow \infty} \lim_{t \rightarrow 1} (1-t) \Phi_{\mu_n, v_n}(t, \varrho) = \lim_{t \rightarrow 1} (1-t) \Phi_{\mu, v}(t, \varrho),$$

if  $\lim_{n \rightarrow \infty} \mu_n = \mu$  and  $\lim_{n \rightarrow \infty} v_n = v$ .

As a first step, we will prove  $\mathcal{L}_{\pi_n}(t, 0) - \log \sqrt{1-t}$  converges to  $\mathcal{L}_{\pi}(t, 0) - \log \sqrt{1-t}$  for sequences  $(\pi_n)_{n=1}^{\infty}$  of probability distributions with  $\lim_{n \rightarrow \infty} \pi_n = \pi$ .

PROPOSITION 1. Let  $(\pi_n)_{n=1}^{\infty}$  and  $\pi$  be distribution functions such that  $\int x d\pi_n = \int x d\pi = 0$  and  $\int x^2 d\pi_n = \sigma_n^2 < \infty$ ,  $\int x^2 d\pi = \sigma^2 < \infty$ . Assume  $\lim_{n \rightarrow \infty} \pi_n = \pi$  (weakly), then for all  $|t| < 1$

$$(12) \quad \lim_{n \rightarrow \infty} \log \sqrt{1-t} - \mathcal{L}_{\pi_n}^{\pm}(t, 0) = \log \sqrt{1-t} - \mathcal{L}_{\pi}^{\pm}(t, 0) \quad \text{and} \\ \lim_{n \rightarrow \infty} \lim_{t \rightarrow 1} (\log \sqrt{1-t} - \mathcal{L}_{\pi_n}^{\pm}(t, 0)) = \lim_{t \rightarrow 1} (\log \sqrt{1-t} - \mathcal{L}_{\pi}^{\pm}(t, 0)).$$

PROOF. We will show the proposition for  $\mathcal{L}_{\pi_n}^{-}$  only. Observe first

$$\mathcal{L}_{\pi_n}^{-}(t, 0) - \log \sqrt{1-t} = \sum_{m=1}^{\infty} \frac{t^m}{m} \left( \mathbb{P} \left( \sum_{i=1}^m X_i^{(n)} < 0 \right) - \frac{1}{2} \right),$$

where  $\mathcal{L}(X_i^{(n)}) = \pi_n$ .

Spitzer (1976) proved that  $\sum_{m=1}^{\infty} \frac{t^m}{m} \left( \mathbb{P} \left( \sum_{i=1}^m X_i < 0 \right) - \frac{1}{2} \right)$  converges (as  $t \rightarrow 1$ ) to

$$\sum_{m=1}^{\infty} \frac{1}{m} \left( \mathbb{P} \left( \sum_{i=1}^m X_i < 0 \right) - \frac{1}{2} \right) \quad \text{and} \quad \sum_{m=1}^{\infty} \frac{1}{m} \left( \mathbb{P} \left( \sum_{i=1}^m X_i < 0 \right) - \frac{1}{2} \right) < \infty$$

provided  $E[X_1] = 0$  and  $E[X_1^2] = \sigma^2 < \infty$  (P4 and P5 Chapter IV, Section 18), see also Rosén's (1962) proof for the absolute convergence of the above series.

Because of continuity of convolutions it is easy to see that (12) holds for all  $|t| < 1$ , even for the cases  $\mu\{0\} > 0$  or  $\nu\{0\} > 0$ . As the limiting function  $\mathcal{L}_{\pi}^{-}(t, 0) - \log \sqrt{1-t}$  and the functions  $\mathcal{L}_{\pi_n}^{-}(t, 0) - \log \sqrt{1-t}$  are continuous for all  $t \in \{t: |t| < 1 \text{ or } t = 1\}$  (P4 of Spitzer (1976)) the claim follows. ■

A similar result holds for the functions  $\mathcal{L}_{\pi_n}^{+}(t, 0)$ .

PROPOSITION 2. Under the conditions of Proposition 1, we have

$$\lim_{n \rightarrow \infty} (1-t)^{-1/2} \exp(\mathcal{L}_{\pi_n}^{\pm}(t, 0)) = (1-t)^{-1/2} \exp(\mathcal{L}_{\pi}^{\pm}(t, 0))$$

for all  $q \in (0, \infty)$ .

PROOF. We have for  $|t| < 1$ :

$$\begin{aligned} & \mathcal{L}_{\pi_n}^{+}(t, 0) - \frac{1}{2} \log(1-t) = \\ &= \sum_{m=1}^{\infty} \frac{t^m}{m} \left( \mathbb{P} \left( \sum_{i=1}^m X_i^{(n)} \geq 0 \right) - \frac{1}{2} \right) + \sum_{m=1}^{\infty} \frac{t^m}{2m} (q^{(1-t)} - 1), \end{aligned}$$

here  $\mathcal{L}(X_i^{(n)}) = \pi_n$ . Hence

$$\lim_{t \rightarrow 1} \mathcal{L}_{\varrho^{1-t}\pi_n}^+(t, 0) - \frac{1}{2} \log(1-t) = \sum \frac{1}{m} \left( \mathbb{P}(\sum X_i^{(n)} \geq 0) - \frac{1}{2} \right) + \frac{1}{2} (1 - \log \varrho).$$

Therefore, Proposition 2 follows from Proposition 1.

We are left to show:

$$\lim_{n \rightarrow \infty} (1-t) \lambda_{\mu_n, v_n}^-(t, \varrho, 0) = (1-t) \lambda_{\mu, v}^-(t, \varrho, 0)$$

for all  $|t| < 1$  and  $t \rightarrow 1$ ,  $\mu_n \rightarrow \mu$  and  $v_n \rightarrow v$ . As  $\lambda_{\mu, v}^-(t, \varrho, s) = \int_{-\infty}^{0^-} e^{sx} dA_{\mu, v}(t, \varrho, x)$ , where  $A_{\mu, v}(t, \varrho, \cdot)$  is the distribution function of a finite (for  $|t| < 1$ ) measure having Laplace transform  $\lambda_{\mu, v}(t, \varrho, s)$ , it is sufficient to show

$$(13) \quad \lim_{n \rightarrow \infty} (1-t) \lambda_{\mu_n, v_n}(t, \varrho, s) = (1-t) \lambda_{\mu, v}(t, \varrho, s),$$

for all  $|t| < 1$  and  $t \rightarrow 1$ , and all imaginary values of  $s$ . This will be demonstrated in the next Proposition.

PROPOSITION 3. Let  $(\mu_n)_1^\infty$ ,  $\mu$ ,  $(v_n)_1^\infty$  and  $v$  be distribution functions such that

$$\int x d\mu_n = \int x d\mu = \int x dv_n = \int x dv = 0, \quad \int x^2 d\mu_n < \infty, \\ \int x^2 d\mu < \infty, \quad \int x^2 dv_n < \infty \quad \text{and} \quad \int x^2 dv < \infty.$$

Assume further  $\lim \mu_n = \mu$  and  $\lim v_n = v$  (weakly), then for all  $|t| < 1$ :

$$\lim_{n \rightarrow \infty} (1-t) \lambda_{\mu_n, v_n}(t, \varrho, s) = (1-t) \lambda_{\mu, v}(t, \varrho, s)$$

and

$$\lim_{n \rightarrow \infty} \lim_{t \rightarrow 1} (1-t) \lambda_{\mu_n, v_n}(t, \varrho, s) = \lim_{t \rightarrow 1} (1-t) \lambda_{\mu, v}(t, \varrho, s)$$

if any of the following three conditions holds:

- (a)  $\mu_n$ ,  $\mu$ ,  $v_n$  and  $v$  are non-lattice distributions,
- (b)  $\mu$  or  $v$  is lattice but  $\{r: \mu(r) \neq 0\} \cap \{r: v(r) \neq 0\} = \emptyset$ ,
- (c)  $\mu_n$ ,  $\mu$ ,  $v_n$  and  $v$  are lattice with  $\hat{\mu}_n(s) = 1 \Leftrightarrow \hat{\mu}(s) = 1$ , and  $\hat{v}_n(s) = 1 \Leftrightarrow \hat{v}(s) = 1$ .

(For any distribution function  $\pi$ , we denote by  $\hat{\pi}$  its two-sided Laplace transform).

PROOF. Consider first the following expressions for  $|t| < 1$ :

$$\lim_{n \rightarrow \infty} (1-t)^{-1/2} \exp \mathcal{L}_{\mu_n}^+(t, s)$$

and

$$\lim_{n \rightarrow \infty} (1-t)^{-1/2} \exp \mathcal{L}_{\varrho^{1-t}v_n}^-(t, s) \quad (s \text{ imaginary}).$$

Similarly to the situation in Proposition 1, it follows from the continuity of the convolution that

$$\lim_{n \rightarrow \infty} (1-t)^{-1/2} \exp \mathcal{L}_{\mu_n}^+(t, s) = (1-t)^{-1/2} \exp \mathcal{L}_{\mu}^+(t, s)$$



and

$$\lim_{n \rightarrow \infty} (1-t)^{-1/2} \exp \mathcal{L}_{\hat{\pi}_n}^-(t, s) = (1-t)^{-1/2} \exp \mathcal{L}_{\pi}^-(t, s).$$

Let  $\pi \in \{\mu, \mu_n: i=1, 2, \dots\}$  and let  $s$  be such that  $\hat{\pi}_n(s) \neq 1$  and  $\hat{\pi}(s) \neq 1$ , then ( $|t| < 1$ )

$$\begin{aligned} \mathcal{L}_{\pi}^+(t, s) - \frac{1}{2} \log(1-t) &= - \sum_{m=1}^{\infty} \frac{t^m}{m} \int_0^{\infty} (1 - \cos(sx)) d\pi^{*m}(x) + \\ &+ \sum_{m=1}^{\infty} \frac{t^m}{m} \left( P\left(\sum_{i=1}^m X_i \equiv 0\right) - \frac{1}{2} \right) + i \sum_{m=1}^{\infty} \frac{t^m}{m} \int_0^{\infty} \sin(sx) d\pi^{*m}(x), \end{aligned}$$

where  $\mathcal{L}(X_i) = \pi$  and  $\mathcal{L}(\sum_{i=1}^m X_i) = \pi^{*m}$ ; similarly

$$\begin{aligned} \mathcal{L}_{\pi}^-(t, s) - \frac{1}{2} \log(1-t) &= - \sum_{m=1}^{\infty} \frac{t^m}{m} \int_{-\infty}^{0^-} (1 - \cos(sx)) d\pi^{*m}(x) + \\ &+ \sum_{m=1}^{\infty} \frac{t^m}{m} \left( P\left(\sum_{i=1}^m X_i < 0\right) - \frac{1}{2} \right) + i \sum_{m=1}^{\infty} \frac{t^m}{m} \int_{-\infty}^{0^-} \sin(sx) d\pi^{*m}(x). \end{aligned}$$

We have:

$$(14) \quad \mathcal{L}_{\pi}^+(t, s) + \mathcal{L}_{\pi}^-(t, s) - \log(1-t) = \sum_{m=1}^{\infty} \frac{t^m}{m} (1 - \operatorname{Re} \hat{\pi}^m(s)) \rightarrow \infty \quad \text{for } t \rightarrow 1.$$

Assume now

$$\lim_{t \rightarrow 1} \sum_{m=1}^{\infty} \frac{t^m}{m} \int_0^{\infty} (1 - \cos(sx)) d\pi^{*m}(x) < \infty.$$

Then, we have from (14)

$$\lim_{t \rightarrow 1} \sum_{m=1}^{\infty} \frac{t^m}{m} \int_{-\infty}^{0^-} (1 - \cos(sx)) d\pi^{*m}(x) \Big/ \sum_{m=1}^{\infty} \frac{t^m}{m} = 1.$$

But

$$\int_{-\infty}^{0^-} (1 - \cos(sx)) d\pi^{*m}(x) \leq P\left(\sum_{i=1}^m X_i < 0\right)$$

and

$$\lim_{m \rightarrow \infty} P\left(\sum_{i=1}^m X_i < 0\right) = \frac{1}{2},$$

hence a contradiction.

Similarly one brings the assumption  $\lim_{t \rightarrow 1} \sum_{m=1}^{\infty} \frac{t^m}{m} \int_{-\infty}^{0^-} (1 - \cos(sx)) d\pi^{*m}(x) < \infty$  to a contradiction. Thus we have proved

$$\lim_{t \rightarrow 1} (1-t)^{-1/2} \exp \mathcal{L}_{\mu_n}^+(t, s) = \lim_{t \rightarrow 1} (1-t)^{-1/2} \exp \mathcal{L}_{\mu}^+(t, s) = 0.$$

Similarly one proves

$$\lim_{t \rightarrow 1} (1-t)^{-1/2} \exp \mathcal{L}_{\varrho^{1-t} \nu_n}(t, s) = \lim_{t \rightarrow 1} (1-t)^{-1/2} \exp \mathcal{L}_{\varrho^{1-t} \nu}(t, s) = 0$$

(see also the proof of Proposition 2).

We are left to prove the claim of Proposition 3 for points with  $\mu_n(s)=1$  and  $\mu(s)=1$  or  $\nu_n(s)=1$  and  $\nu(s)=1$ . Here we have, for any distribution with  $\hat{\pi}(s)=1$

$$|1 = |\hat{\pi}(s)| \leq |\hat{\pi}^+(s)| + |\hat{\pi}^-(s)| \leq P(X \geq 0) + P(X < 0) = 1$$

$$(\mathcal{L}(X) = \pi) \text{ and hence } \hat{\pi}^+(s) = P(X \geq 0),$$

the claim of the proposition follows then from Propositions 1 and 2.

We showed now the validity of the arcsine-law for the oscillating random walk if the underlying distributions are either non-lattice or live on disjoint lattices. We are left to consider the lattice case, which will be done in the next section.

## 5. The lattice case

Here we consider random variables  $Y_{1,1}$  and  $Y_{2,1}$  for which there exists a common lattice. As we are only interested in the random variable  $N_n$ , as defined in (2), we may as well restrict ourselves to random variables with values in  $\mathbb{Z}$ , the set of integers.

Similar to the treatment in the previous section, first we restrict ourselves to increment distributions  $\mu$  and  $\nu$ , for which we can obtain explicit formulae for  $Q^+(t, \varrho, s)$  and  $Q^-(t, \varrho, s)$ . More precisely we assume

$$\hat{\nu}(s) = \nu^*(e^s) = E[\exp(sY_{1,1})] = \sum_{j=-s^+}^{s^+} \nu(j) e^{sj}, \quad (15)$$

$$\hat{\mu}(s) = \mu^*(e^s) = E[\exp(sY_{2,1})] = \sum_{j=-s^-}^{s^-} \mu(j) e^{sj}.$$

We then find (compare Kemperman (1974))

$$\exp \mathcal{L}_{\mu}^-(t, s) = \sum_{v=-1}^{-s^-} (e^s / (e^s - \xi_v(t)))$$

$$\exp \mathcal{L}_{\mu}^+(t, s) = (-t\mu(r^-))^{-1} \prod_{v=1}^{s^-} (e^s - \xi_v(t))^{-1}$$

$$\exp(\mathcal{L}_{\nu}^-(t, s)) = \prod_{v=-1}^{-s^-} (e^s / (e^s - \eta_v(\varrho t)))$$

$$\exp(\mathcal{L}_{\nu}^+(t, s)) = (-\varrho t\nu(r^+))^{-1} \prod_{v=1}^{s^+} (e^s - \eta_v(\varrho t))^{-1}, \quad (16)$$

where  $\xi_v(t)$  are the  $(r^- + s^-)$  roots of the equation  $\mu^*(\xi_v) = 1/t$ , where it is assumed  $|\xi_i| \leq |\xi_j|$  if  $i \leq j$  ( $|t| < 1$ ); similarly the  $\eta_v(t)$  are the  $(r^+ + s^+)$  roots of the equation  $\nu^*(\eta_v) = 1/t$ , again  $|\eta_i| \leq |\eta_j|$  for  $i \leq j$  ( $|t| < 1$ ) is assumed.

From this we obtain:

$$(17) \quad (\delta_z \lambda)^- = \sum_{k=-1}^{-s^+} \frac{\eta_k^{z+1}(\varrho t)}{e^s - \eta_k(\varrho t)} \prod_{v=1}^{r^-} (\xi_v(t) / (\xi_v(t) - \eta_k(\varrho t))) \times \\ \times \prod_{\substack{v=-1 \\ v \neq k}}^{-s^+} (\eta_k(\varrho t) / (\eta_k(\varrho t) - \eta_v(\varrho t))),$$

and

$$(\delta_z \lambda)^+ = (\delta_z \lambda) - (\delta_z \lambda)^-$$

if  $z \geq 0$ , for  $z < 0$  we have

$$(18) \quad (\delta_z \lambda)^+ = \sum_{k=1}^{r^-} \frac{\xi_k^{z+1}(t)}{\xi_k(t) - e^s} \prod_{\substack{v=1 \\ v \neq k}}^{r^-} (\xi_v(t) / (\xi_v(t) - \xi_k(t))) \times \\ \times \prod_{v=-1}^{-s^+} (\xi_k(t) / (\xi_k(t) - \eta_v(\varrho t))) \\ (\delta_z \lambda)^- = (\delta_z \lambda) - (\delta_z \lambda)^+.$$

Formulae (16) together with (17) or (18) depending on  $z \geq 0$  or  $z < 0$  give an explicit formulae for  $\Phi(t, \varrho)$ , see (5). Using again Corollary 1 of Appendix A, we obtain the same limiting distribution for  $N_n/n$  as in the non-lattice case.

Distributions fulfilling (15) are arbitrary distributions on  $\mathbf{Z}$  with finite support, which clearly approximate any distribution on  $\mathbf{Z}$ . Hence we can apply Proposition 1, 2 and 3 of the previous section and we proved the following theorem:

**THEOREM 2.** Let  $(Y_{1,i})_{i=1}^\infty$  and  $(Y_{2,i})_{i=1}^\infty$  be two independent sequences of independent identically distributed random variables with  $E[Y_{1,1}] = E[Y_{2,1}] = 0$ ,  $E[Y_{1,1}^2] = \sigma_1^2$  and  $E[Y_{2,1}^2] = \sigma_2^2$ , let  $N_n$  be defined by (2), then  $\lim_{n \rightarrow \infty} \mathcal{L}(N_n/n) = \mathcal{L}(Y)$ , where  $Y$  is a random variable with density  $g(y)$  (with respect to the Lebesgue measure) and

$$g(y) = \pi^{-1} R(y(1-y))^{-1/2} (y + (1-y)R^2)^{-1} I_{(0,1)}(y).$$

## Appendix A

A Tauberian theorem for a special class of analytic functions is derived:

**PROPOSITION 1.** Let  $f(t, z) = \sum_{m=0}^\infty a_m(z) t^m$  be an analytic function in  $z$  and  $t$ , for  $|1-z| < 1$  and  $|t| < 1$ . Assume  $a_m(z) = \sum_{i=1}^\infty b_i^{(m)} (1-z)^i$  is analytic for  $|z| < 1$ . Then we can write

$$f(t, z) = \sum_{i=0}^\infty \left( \sum_{m=0}^\infty b_i^{(m)} t^m \right) (1-z)^i \quad (|t| < 1, |z| < 1).$$

Assume further

- (a)  $b_i^{(m)} \geq 0$ ,  $i \in \mathbb{N}$ ,  $m \in \mathbb{N}$ ,  
 (b)  $(b_i^{(m)})_{m=0}^\infty$  is a monotonely nondecreasing sequence,  
 (c)  $\lim_{t \rightarrow 1} (1-t)f(t, 1-\omega(1-t)) = L_\omega$  for all  $|\omega|=1$ .

Then

$$\lim_{m \rightarrow \infty} \sum_{i=0}^{\infty} (i!) b_i^{(m)} m^{-i} = L_1.$$

An important corollary to this proposition connects limits of generating functions of sequences of random variables  $(X_i)_i^\infty$  with the limits of the Stieltjes transforms of the variables  $(X_i/i)_i^\infty$ .

COROLLARY 1. Assume  $(X_i)_i^\infty$  is a stochastically increasing sequence of non-negative random variables, and let

$$g(t, z) = \sum_{i=0}^{\infty} E[\varrho^{(1-z)X_i}] t^i,$$

then, for all  $\varrho$  such that  $X_i \log \varrho < i$  (a.s.):

$$\lim_{i \rightarrow \infty} E[(1 - (\log \varrho) X_i / i)^{-1}] = \lim_{t \rightarrow 1} (1-t)g(t, t)$$

provided  $\lim_{t \rightarrow 1} (1-t)g(t, 1-\omega(1-t))$  exists for all  $|\omega|=1$ .

We first prove the proposition: Consider

$$f(t, z) = \sum_{i=0}^{\infty} (1-t)^i \sum_{m=0}^{\infty} b_i^{(m)} z^m,$$

then

$$\lim_{t \rightarrow 1} (1-t)f(t, 1-\omega(1-t)) = L_\omega, \quad (|\omega|=1)$$

implies  $(b_i^{(m)} \geq 0)$

$$\lim_{t \rightarrow 1} (1-t)^{i+1} \sum_{m=0}^{\infty} b_i^{(m)} t^m = L'_i$$

for every  $i \in \mathbb{N}$ , further  $\sum_{i=0}^{\infty} L'_i = L_1$ .

The proof of the above claim lies essentially in the continuity theorem for Fourier series, (see e.g. Feller (1971) Theorem 2, p. 431). Using a Tauberian theorem (Feller (1971) Theorem 5, p. 447) we obtain  $\lim_{m \rightarrow \infty} i! m^{-i} b_i^{(m)} = L'_i$ , as

$$\lim_{m \rightarrow \infty} \sum_{i=0}^{\infty} i! (m)^{-i} b_i^{(m)} = \sum_{i=0}^{\infty} \lim_{m \rightarrow \infty} i! (m)^{-i} b_i^{(m)}$$

(positivity of all summands), the theorem is proved.

To prove the corollary one only has to check the equality  $b_i^{(m)} = E[X_m^i] (\log \varrho)^i / i!$ .

## Appendix B

The denseness of the family  $\mathcal{F}$  of positive finite measures over  $\mathbf{R}^+$  with densities  $\sum_{i=0}^m m_i e^{-\lambda_i x}$  (with respect to the Lebesgue measure) is demonstrated.

PROPOSITION. Let  $\mathcal{F}$  be the family of positive finite measures over  $\mathbf{R}^+$  with densities  $\sum_{i=0}^m m_i e^{-\lambda_i x}$  (density with respect to the Lebesgue measure), then  $\mathcal{F}$  is dense in the space of all finite measures over  $\mathbf{R}^+$ .

PROOF. The algebra generated by  $\{1, e^{-x}\}$  is dense in the space of all measurable functions over any compact set contained in  $\mathbf{R}^+$  in the supremum topology (Stone—Weierstraß). Therefore  $\mathcal{F}$  is vag dense in the space of finite measures over  $\mathbf{R}^+$ . Then Theorem 45,7 on p. 188 of Bauer (1968) proves the proposition.

## REFERENCES

- [1] BAUER, H., *Wahrscheinlichkeitstheorie und Grundzüge der Maßtheorie*, de Gruyter, Berlin, 1968. MR 39#983.
- [2] FELLER, W., *An introduction to probability theory and its application II*, J. Wiley, New York, 1971. MR 42#5292.
- [3] KEILSON, J. and WELLNER, J. A., Oscillating Brownian motion, *J. Appl. Probability* 15 (1978), 300—310. MR 57#14164.
- [4] KEMPERMAN, J. H. B., *The passage problem for a stationary Markov chain*, The University of Chicago Press, Chicago, 1961. MR 22#9992.
- [5] KEMPERMAN, J. H. B., The oscillating random walk, *Stochastic Processes Appl.* 2 (1974), 1—29. MR 50#14940.
- [6] LAMPERTI, J., An occupation time theorem for a class of stochastic processes, *Trans. Amer. Math. Soc.* 88 (1958), 380—387. MR 20#1372.
- [7] REIMNITZ, P., A study of Two Armed Bandits, *unpublished Ph. D.-thesis*, University of Rochester, Rochester, N. Y. (1977).
- [8] ROSÉN, B., On the asymptotic distribution of sums of independent identically distributed random variables. *Ark. Mat.* 4 (1962), 323—332. MR 25#5528.
- [9] SPITZER, F. L., *Principles of Random Walk*, Springer-Verlag, New York, Heidelberg, Berlin, 1976. MR 52#9383.
- [10] WIDDER, D. V., *The Laplace Transform*, Princeton University Press, Princeton, N. Y., 1972.

(Received March 7, 1984)

INSTITUT FÜR MEDIZINISCHE STATISTIK  
 DOKUMENTATION UND DATENVERARBEITUNG  
 DER UNIVERSITÄT BONN  
 D-5300 BONN 1  
 SIGMUND-FREUD-STRASSE 25  
 FEDERAL REPUBLIC OF GERMANY



# ON THE SPECTRUM OF A CLASS OF INTEGRAL TRANSFORMS II

B. P. DUGGAL

## Abstract

Let  $C(L^p)$ ,  $1 < p < \infty$ , denote the class of linear transformations which are continuous on  $L^p$  into itself. The mapping  $T \in C(L^p)$  is said to belong to the class  $G_0$  if  $T$  satisfies the functional relation  $Tt(a) = m(a)t(a)T$ , where  $t(a)$  denotes the operator of dilatation by amount  $a$  ( $a \neq 0$ ), and where  $m(a) = (\text{sgn } a)$  or 1. It is shown that if  $T^2 = bI$  for some scalar  $b$ , then the spectrum  $\sigma(T)$  of  $T$  consists purely of the point spectrum  $\sigma_p(T)$ . Furthermore,  $\sigma_p(T) = \{\pm \sqrt[2]{b}\}$ . If  $T' = cT$  for some scalar  $c$  such that  $|c| = 1$ , and if  $T$  is invertible, then it is shown that either  $\sigma(T) = \sigma_p(T)$  or  $\sigma(T) = \sigma_c(T)$  (=the continuous spectrum of  $T$ ).

## 1. Introduction

Let  $L^p = L^p(-\infty, \infty)$ , and let  $C(L^p)$  denote the class of linear transformations which are continuous on  $L^p$  into itself. We say that the linear transformation  $T \in C(L^p)$ ,  $1 < p < \infty$ , belongs to the class  $G_0$  if  $T$  satisfies the functional equation

$$(1) \quad Tt(a) = m(a)t(a)T, \quad \text{all real } a \neq 0,$$

where  $t(a)$  denotes the operator of dilatation by amount  $a$ , and where  $m(a) = (\text{sgn } a)$  or 1. It is known (see [2]) that if  $T \in G_0$ , then there exists a Lebesgue measurable function  $K$  such that

$$(2) \quad \int_0^u Tf(t) dt = \int_{-\infty}^{\infty} m(x)K(ux^{-1})f(x) dx.$$

A number of the important (in applications) integral transforms belong to the class  $G_0$ . Thus, for example, the Hilbert transform  $H$ , the Stieltjes transform  $S$ , and the modified Hilbert transform  $H^{(\alpha)}$  defined (respectively) by

$$Hf(u) = (1/\pi)(\text{P.V.}) \int_{-\infty}^{\infty} (u-x)^{-1}f(x) dx;$$

$$Sf(u) = (1/\pi) \int_0^{\infty} (u+x)^{-1}f(x) dx;$$

$$H^{(\alpha)}f(u) = (1/\pi)|u|^{-\alpha}(\text{P.V.}) \int_{-\infty}^{\infty} |x|^{\alpha}(u-x)^{-1}f(x) dx, \quad -1/p < \alpha < 1/p',$$

all belong to  $G_0$ .

1980 *Mathematics Subject Classification*. Primary 44A05; Secondary 47A10.  
Key words and phrases. Integral transform, adjoint,  $L^p$ -space, spectrum.

The study of the spectrum  $\sigma(T)$  of particular mappings  $T \in G_0$  has been carried out by a number of authors, amongst them Pollard [9], de Snoo [11], and Juberg [7]. Starting from the view point of the functional equation characterizing mappings  $T \in G_0$ , a study of the spectrum of the class  $G_0$  was initiated by the author in [4]. We showed that  $\sigma(T) = \sigma_p(T) \cup \sigma_c(T)$ , where  $\sigma_p(T)$  and  $\sigma_c(T)$  denote the point spectrum and the continuous spectrum of  $T$ , and that if  $T^2 = bI$  on  $L^p$ , then  $\sigma_p(T)$  is not empty and  $\sigma(T) \subseteq \{\pm\sqrt{b}\}$ . We went onto state (see Theorem 3) that if  $T' = -T$ , then there exists a scalar  $b$  such that  $\sigma(T) = \sigma_p(T) = \{\pm\sqrt{b}\}$ . Unfortunately, there is an error in the proof here. In this paper we locate the error, and give a number of sufficient conditions for the theorem (*loc. cit.*) to hold. Indeed, we show that if  $T \in G_0$  is 0-adjoint (i.e.  $T' = cT$  for some scalar  $c$  such that  $|c| = 1$ ), and invertible, then either  $\sigma(T) = \sigma_p(T)$  or  $\sigma(T) = \sigma_c(T)$  (exclusive 'or'). Also, we strengthen Corollary following Theorem 3 of [4] to show that if  $T^2 = bI$  for some scalar  $b$ , and if  $T$  is not scalar type, then  $\sigma(T) = \sigma_p(T) = \{\pm\sqrt{b}\}$ . A number of examples are considered to illustrate the results.

## 2. Notation and complementary results

In addition to the notation already introduced, the following further notation will be used. The empty set will be denoted by  $\varnothing$  (as apart from the function  $\varphi$ ), and the set of real (complex) numbers will be denoted by  $\mathbf{R}$  (resp.,  $\mathbf{C}$ ). The resolvent set of  $T$  will be denoted by  $\varrho(T)$ . We assume, once for all, that  $1 < p < \infty$  and that  $T \in G_0$ . Also,  $T$  will be assumed to have the integral representation (2) with  $m(x) = 1$ . (There is no loss of generality in assuming  $m(x) = 1$ : our results hold just as well for the case in which  $m(x) = (\text{sgn } x)$ .) The identity map will be denoted by  $I$ , the mapping (Banach space) adjoint to  $T$  will be denoted by  $T'$ , and the index conjugate to  $p$  will be denoted by  $p'$ . Henceforth, whenever the limits extend over all of  $\mathbf{R}$ , the limits will be omitted from the integrals. We note that all our equalities involving functions are to be considered as holding a.e. only.

The mapping  $T$  is said to be 0-adjoint if there exists a scalar  $c$ ,  $|c| = 1$ , such that  $T' = cT$  (see [2]). We assume, henceforth, that  $c = 1$ : this does not involve any loss of generality (for if  $c \neq 1$ , then we define  $P = dT$  with  $\bar{d}c = d$ , and consider  $P = dT = \bar{d}cT = \bar{d}T' = P'$ ). Let  $T$  be invertible with  $T^{-1} = S$ . Then, this is easily verified,  $S$  satisfies a functional equation of type (1), and so  $S \in G_0$ . Hence  $S$  has an integral representation of the type

$$(3) \quad \int_0^u Sf(t) dt = \int Q(ux^{-1})f(x) dx, \quad u \in \mathbf{R},$$

for some Lebesgue measurable function  $Q$ . Let  $e$  be the function  $e(x) = 1$  if  $0 < x < 1$ , and zero otherwise. It is not difficult to verify (see [2] and [3]) that if

$$(4) \quad \int_0^u \psi(t^{-1}) dt = u^{-1} \int_0^{1/u} K(t^{-1}) dt;$$

$$(5) \quad \int_0^u \varphi(t^{-1}) dt = u^{-1} \int_0^{1/u} Q(t^{-1}) dt,$$



then  $Te(t)=\psi(t^{-1})$  and  $Se(t)=\varphi(t^{-1})$ . The mappings  $T'$  and  $S'$  then have the integral representation

$$(6) \quad \int_0^{\infty} T'f(t) dt = \int \psi(ux^{-1})f(x) dx, \quad u \in \mathbf{R};$$

$$(7) \quad \int_0^{\infty} S'f(t) dt = \int \varphi(ux^{-1})f(x) dx, \quad u \in \mathbf{R}.$$

### 3. An example

During the course of the proof of Theorem 3 (Theorem 4) of [4] we stated that by Theorem 3.2 of [3] if  $T$  (resp.,  $T'T$ ) is an invertible 0-adjoint mapping belonging to  $G_0$ , then there exists a non-zero scalar  $b$  such that  $T^2=bI$  (resp.,  $(T'T)^2=bI$ ) on  $L^p$ . This is not true, as the following example shows.

Let  $T=I-P$ , where  $P=(2 \cos \pi v/\pi)M_{(v)}^2$  and  $M_{(v)}$  is the Meijer transform

$$M_{(v)}f(x) = \sqrt{\frac{2}{\pi}} \int_0^{\infty} \sqrt{xt} K_v(xt)f(t) dt, \quad |\operatorname{Re} v| < 1.$$

If  $\operatorname{Im} v=0$ ,  $\operatorname{Re} v \neq 0$  and  $v \neq \pm \frac{1}{2}$ , then  $T$  is an 0-adjoint invertible mapping in  $G_0 \cap C(L^2(0, \infty))$ . However,  $T^2 \neq bI$ .

Thus Theorem 3 (Theorem 4) of [4] is valid only for the case in which  $T^2=bI$  (resp.,  $(T'T)^2=bI$ ). The fault here lies with Theorem 3.3 of [3]; more precisely our conclusion there (with  $m(x)=m'(x)=1$ ) that

$$(8) \quad \int \psi(x^{-1})Q(ux^{-1}) dx = \int \varphi(x^{-1})K(ux^{-1}) dx$$

implies

$$(9) \quad \psi(x^{-1})Q(ux^{-1}) = \varphi(x^{-1})K(ux^{-1})$$

(see p. 282) is false. Some additional hypotheses are required. In this section we provide these hypotheses.

Let  $f_v(x)$  denote the function  $f_v(x)=f(v^{-1}x)$ .

**THEOREM 1.** *Let  $T$  be invertible with inverse  $S$ . Then  $T^2=bI$  on  $L^p$  if and only if there exists an  $f \in L^p$  such that  $Tf_v(x)=b^{-1}Sf_v(x)=e_v(x)$ .*

**PROOF.** Let  $T$  and  $S$  have the integral representations (2) and (3). Then the argument of the proof of Theorem 3.3 of [3] (irrespective of whether  $T$  is 0-adjoint or not) shows that (8) holds. A simple change of variable now leads us to

$$\int \psi(vx^{-1})Q(ux^{-1}) dx = \int \varphi(vx^{-1})K(ux^{-1}) dx$$

for all non-zero real  $v$ . As stated in Section 2,  $\psi(x^{-1})=Te(x)$  and  $\varphi(x^{-1})=Se(x)$ , and so, by the functional equation satisfied by  $S$  and  $T$ , we have that  $\psi(vx^{-1})=Te_v(x)$  and  $\varphi(vx^{-1})=Se_v(x)$ . The absolute existence of the integrals in (2) and (3)

implies that both  $K(ux^{-1})$  and  $Q(ux^{-1})$  are in  $L^p$  for each real  $u$ ; hence we have from an application of Parseval's relation that

$$\begin{aligned}\int \psi(vx^{-1})Q(ux^{-1})dx &= \int Te_v(x)Q(ux^{-1})dx = \\ &= \int e_v(t)[T'Q(ux^{-1})](t)dt = \\ &= \int \varphi(vx^{-1})K(ux^{-1})dx = \\ &= \int e_v(t)[S'K(ux^{-1})](t)dt.\end{aligned}$$

Since this holds for all non-zero real  $v$ , we have that

$$(10) \quad [T'Q(ux^{-1})](t) = [S'K(ux^{-1})](t),$$

where the resultant function on either side of the equality is in  $L^p$ . Now let  $f_v(t)$  be an element of  $L^p$  (for each non-zero real  $v$ ). Then we have from (10) and another application of Parseval's relation that

$$\int f_v(t)[T'Q(ux^{-1})](t)dt = \int f_v(t)[S'K(ux^{-1})](t)dt,$$

or

$$\int Q(ux^{-1})Tf_v(x)dx = \int K(ux^{-1})Sf_v(x)dx,$$

i.e.

$$\int_0^v Q(ux^{-1})dx = b \int_0^v K(ux^{-1})dx$$

for all non-zero real  $v$  (and all  $u$ ). Hence  $Q(t) = bK(t)$ , which implies that  $T^2 = bI$ .

To complete the proof, we let  $T^2 = bI$ , and define  $S \in G_0$  by  $S = b^{-1}T$ . Then

$$\int_0^u T^2 e_v(t)dt = \int_0^u [T\psi(vx^{-1})](t)dt = b \int_0^u e_v(t)dt.$$

Hence,  $[T\psi(vx^{-1})](u) = be_v(u)$ . Setting  $b^{-1}\psi(vx^{-1}) = f_v(x)$ , we see that we have found the required function.

**THEOREM 2.** Let  $T$  be 0-adjoint invertible with inverse  $S$ . Set

$$\{Q(x^{-1})K(ux^{-1}) - K(x^{-1})Q(ux^{-1})\} = G(u, x).$$

If  $\int G(u, x)dx = 0$  for a.a.u implies that  $G(u, x) = 0$  a.e., then there exists a scalar  $b$  ( $\neq 0$ ) such that  $T^2 = bI$ .

**PROOF.** Clearly, equality (8) holds. Since  $T$  is 0-adjoint,  $S$  is 0-adjoint, and so we have from (6) and (7) that

$$\int K(ux^{-1})f(x)dx = \int \psi(ux^{-1})f(x)dx;$$

$$\int Q(ux^{-1})f(x)dx = \int \varphi(ux^{-1})f(x)dx,$$

for each  $f \in L^p$  and all  $u$ . Letting  $f(x) = e(x)$ , (4) and (5) now imply that  $\varphi(x^{-1}) = Q(x^{-1})$  and  $\psi(x^{-1}) = K(x^{-1})$ . Substituting in (8), we then have from the hypothesis that

$$K(x^{-1})Q(ux^{-1}) = Q(x^{-1})K(ux^{-1})$$

for all  $u$ . Now set  $ux^{-1}=a$ , a some non-zero real number, and then  $x^{-1}=t$ ; we have

$$Q(t) = (Q(a)/K(a))K(t) = bK(t)$$

for some scalar  $b$  ( $\neq 0$ ). This completes the proof.

**THEOREM 3.** *Let  $T$  be invertible with inverse  $S$ . Then either of the following conditions implies that  $T^2=bI$ .*

$$(i) \quad \int \psi(ux^{-1})K(x^{-1}) = b\{\min(u, 1) - \min(u, 0)\}.$$

(ii) *There exist functions  $\varphi_1(t)=\varphi_1(t)e(a^{-1}t)$  and  $\varphi_2(t)$ , both in  $L^p$ , such that  $T'\varphi_1(t)=\varphi_2(t)$  and  $T'\varphi_2(t)=b\varphi_1(t)$ .*

$$(iii) \quad \int G(u, x) dx = \int \{\varphi(x^{-1})K(ux^{-1}) - \psi(x^{-1})Q(ux^{-1})\} dx = 0$$

for a.a.u implies  $G(u, x)=0$  a.e.

For the proof of cases (i) and (ii) we refer to [2], Corollaries 8 and 9; the proof of (iii) follows from an argument similar to that used in the proof of Theorem 2.

#### 4. Spectrum: case $\sigma_p(T) \neq \emptyset$

Throughout the following we assume that  $T$  is not a scalar multiple of the identity, i.e.  $T \neq aI$  for some  $a \in \mathbb{C}$ .

**THEOREM 4.** *If  $T^2=bI$ , then  $\sigma(T)=\sigma_p(T)=\{\pm\sqrt{b}\}$ .*

**PROOF.** By the Corollary following Theorem 3 of [4],  $\sigma(T) \subseteq \{\pm\sqrt{b}\}$  and  $\sigma_p(T)$  is not empty. (We remark that the corollary *loc. cit.* holds true.) We show that  $\sigma(T)=\{\pm\sqrt{b}\}$ .

Since  $\sigma(T)$  is not empty, at least one of the points  $-\sqrt{b}$  and  $+\sqrt{b}$  is in  $\sigma(T)$ ; say  $-\sqrt{b} \in \sigma(T)$ . Since our assertion is trivially true if  $+\sqrt{b}$  also belongs to  $\sigma(T)$ , we assume that  $+\sqrt{b} \notin \sigma(T)$ . The mapping  $(\sqrt{b}I - T)$  is then (continuously) invertible. Since  $T^2=bI$ ,  $(bI - T^2)=(\sqrt{b}I - T)(\sqrt{b}I + T)=0$ , and hence  $(\sqrt{b}I + T)=0$ . But then  $T$  is a scalar multiple of the identity. Hence  $+\sqrt{b} \in \sigma(T)$ .

Now since  $\sigma_p(T)$  is not empty, at least one of the points  $-\sqrt{b}$  and  $+\sqrt{b}$  is in  $\sigma_p(T)$ . Assume that  $+\sqrt{b} \in \sigma_p(T)$  and that  $-\sqrt{b} \notin \sigma_p(T)$ . (Recall that  $\sigma(T) = \sigma_p(T) \cup \sigma_c(T)$ .) Then the mapping  $(T + \sqrt{b}I)$  is one to one, and the range of  $(T + \sqrt{b}I)$  is dense in  $L^p$ . Hence, to each  $f \in L^p$  there corresponds just one  $g \in L^p$  such that

$$(11) \quad Tf + \sqrt{b}f = g,$$

the set of such  $g$  being dense in  $L^p$ . Applying  $T$  to both sides of (11) we have that  $b f + \sqrt{b} Tf = Tg$ , and so that  $Tg = \sqrt{b}g$ . But then  $T = \sqrt{b}I$  on a dense subset of  $L^p$ , contradicting thereby the hypothesis that  $T$  is not scalar type. Hence  $-\sqrt{b} \in \sigma_p(T)$ . This completes the proof.

Both  $H$  and  $H^{(a)}$  satisfy the hypotheses of Theorem 2 with  $b = -1$ ; it follows from Theorem 4 that the spectrum of  $H$  and  $H^{(a)}$  consists purely of the point spec-



trum, and is the set  $\{\pm i\}$ . Let  $f \in L^p$ , and let  $U$  be the mapping

$$Uf(z) = (1/2\pi i) \int (t-z)^{-1} f(t) dt, \quad z \in \mathbb{C} \setminus \mathbb{R}.$$

Then  $U^+f(x) = \lim_{y \rightarrow 0^+} Uf(x+iy)$  and  $U^-f(x) = \lim_{y \rightarrow 0^-} Uf(x+iy)$  exist a.e., and satisfy

$$U^+ = (1/2i)(-H+i), \quad U^- = (-1/2i)(H+i), \quad U^+ - U^- = I$$

on  $L^p$  (see [1]). Let  $M_p = \{f \in L^p : U^+f = 0\}$  and  $N_p = \{f \in L^p : U^-f = 0\}$ ; then  $L^p = M_p \oplus N_p$ , and each  $f \in L^p$  has a unique representation of the type  $f = f_1 + f_2$ , where  $f_1 \in M_p$  and  $f_2 \in N_p$ . Since mappings of the class  $G_0$  mutually commute (see [2]), we have that  $HTf_1 = THf_1 = iTf_1$  and that  $HTf_2 = THf_2 = -iTf_2$ .

**THEOREM 5.**  $\sigma_p(TH) \cup \sigma_p(T'H) = \sigma_p(H) \sigma_p(T) = \{\pm i\alpha : \alpha \in \sigma_p(T)\}$ .

**PROOF.** Suppose that  $\alpha \in \sigma_p(T)$ ; then there is a non-trivial  $f \in L^p$  such that  $Tf = \alpha f$ . Let  $f = f_1 + f_2$ , where  $f_1 \in M_p$  and  $f_2 \in N_p$ . Then

$$THf = HTf = \alpha Hf = i\alpha(f_1 - f_2),$$

and so, since  $H^2 = -I$ ,

$$-\alpha f = -Tf = THHf = -iTH(f_1 - f_2).$$

Hence,  $THf_1 = i\alpha f_1$  and  $THf_2 = -i\alpha f_2$ . Now if  $f_1 \neq 0$ , then  $i\alpha \in \sigma_p(TH)$ , and so since  $(TH)' = H'T' = -HT' = -T'H$ ,  $-i\alpha \in \sigma_p(T'H)$ . Again, if  $f_2 \neq 0$ , then  $-i\alpha \in \sigma_p(TH)$ , and so also  $i\alpha \in \sigma_p(T'H)$ . Finally, since  $\sigma(H) = \sigma_p(H) = \{\pm i\}$ , we see that  $\sigma_p(TH) \cup \sigma_p(T'H) = \sigma_p(H) \sigma_p(T)$ .

**COROLLARY 1.** If  $T$  is 0-adjoint, then  $\sigma_p(T)$  is symmetric about the origin.

**PROOF.** Since  $T' = T$ , Theorem 5 implies that  $\sigma_p(TH) = \sigma_p(T) \sigma_p(H) = \{\pm i\alpha : \alpha \in \sigma_p(T)\}$ . Hence  $\sigma_p(HTH) = -\sigma_p(T) = \sigma_p(H) \sigma_p(TH) = \{\pm i\alpha : \alpha \in \sigma_p(T)\}$ .

It is immediate from Corollary 1, and the fact that  $\sigma(T) = \sigma_p(T) \cup \sigma_c(T)$ , that if  $T$  is 0-adjoint and if  $\sigma(T)$  lies to one side of the origin, then  $\sigma_p(T) = \emptyset$ . Let  $R$  be the mapping  $Rf(x) = x^{-1}f(x^{-1})$ , and let  $I_1$  denote the class of mappings obtained by defining  $V \in I_1$  if  $V = TR$  for some  $T \in G_0$  (see [2] for more detail about the class  $I_1$ ). Let  $V \in I_1 \cap C(L^2)$ . Then it can be shown that  $\alpha \in \sigma_p(V)$  if and only if  $\alpha^2 \in \sigma_p(V^2)$  (see the proof of Theorem (3.5) of [5]). Let  $T \in C(L^2)$ .

**COROLLARY 2.** If  $\alpha \in \sigma_p(V)$ , then the set  $\{\pm\alpha, \pm i\alpha\} \in \sigma_p(V)$ .

**PROOF.** Since  $R$  is an isometric isomorphism of  $L^2$  onto itself,  $V = TR \in C(L^2)$ . Since  $RT = T'$  (see [2]),  $V^2 = TRTR = TT'$  is 0-adjoint. Hence if  $\alpha^2 \in \sigma_p(V^2)$ , then  $-\alpha^2 \in \sigma_p(V^2)$ . This, by the remark above, establishes the corollary.

An example of a mapping illustrating Corollary 2 is provided by the Fourier transform  $F$ . Recall that  $1 \in \sigma_p(F)$ ; so, by Corollary 2,  $\{\pm 1, \pm i\} \in \sigma_p(F)$ . Indeed, it is known that  $\sigma(F) = \sigma_p(F) = \{\pm 1, \pm i\}$ .

**COROLLARY 3.** If  $T$  is 0-adjoint, and if  $0 \in \sigma(T)$ , then  $(aI - T)^2 = bI$  cannot hold for any  $a, b \in \mathbb{C}$ .

PROOF. For all  $T$ , not necessarily 0-adjoint,  $(aI - T)^2 = bI$  implies that  $\sigma(aI - T) = \sigma_p(aI - T) = \{\pm\sqrt{b}\}$ . Hence, since  $0 \in \sigma_p(T)$ ,  $\sigma(T) = \sigma_p(T)$  is either the set  $\{0, 2a\}$  or the set  $\{-2a, 0\}$ . Now if  $T$  is 0-adjoint, then, by Corollary 1,  $\sigma(T) = \sigma_p(T) = \{0\}$ . Hence  $T = 0$ . (We note here that because of the mutual commutativity of the mappings  $T \in G_0$ , on the Hilbert space  $L^2$  the mappings  $T$  are normal.)

Although  $(aI - T)^2 = bI$  cannot hold for 0-adjoint  $T$  such that  $0 \in \sigma(T)$ , a slight variation bears fruit. Thus, let  $Af(x) = a(\operatorname{sgn} x)f(x) - Tf(x)$  for 0-adjoint  $T$  such that  $0 \in \sigma(T)$ . If  $A^2 = bI$ , then  $\sigma(A) = \sigma_p(A) = \{\pm\sqrt{b}\}$ , and so  $\sigma(T) = \sigma_p(T) = \{\pm(a - \sqrt{b}), \pm(a + \sqrt{b})\}$ . Since  $0 \in \sigma(T)$ , we have that  $\sigma(T) = \sigma_p(T) = \{0, \pm 2a\}$ . An example of a mapping  $T$  satisfying these conditions is provided by

$$Tf(x) = (2/\pi^2) \int (\log t - \log x)/(t - x)f(t) dt.$$

### 5. Spectrum: case $\sigma_c(T) \neq \varphi$

Once again we assume that  $T$  is not a scalar multiple of the identity. We start by considering the case in which  $T$  is invertible.

THEOREM 6. Let  $T$  be invertible. Then  $\sigma_c(T) \neq \varphi$  if and only if there is a  $V \in G_0$  such that  $0 \in \sigma_c(V)$  and  $T = aI - V$  for some  $a \in \mathbb{C}$ .

PROOF. Suppose that  $\sigma_c(T) \neq \varphi$ . Then there is non-zero scalar  $b$  such that  $b \in \sigma_c(T)$ . Let  $S = T^{-1}$ , and let  $T$  and  $S$  have the integral representations (2) and (3), respectively. Then (8) is satisfied. Now  $b^{-1} \in \sigma_c(S)$ , and so to each  $h \in L^p$  there corresponds just one  $g \in L^p$  such that  $Sh - b^{-1}h = g$ . Letting  $h(x) = \psi(x^{-1})$  ( $\in L^p$ ), we have that

$$\int_0^u [S\psi(x^{-1})](t) dt - b^{-1} \int_0^u \psi(t^{-1}) dt = \int_0^u g(t) dt.$$

Since  $\psi(x^{-1}) = Te(x)$ , this implies that

$$b^{-1} \int_0^u \psi(t^{-1}) dt = \int_0^u \{STe(t) - g(t)\} dt = \int_0^u \{e(t) - g(t)\} dt$$

for all  $u$ . Hence

$$b^{-1}\psi(t^{-1}) = e(t) - g(t).$$

Substituting in (6), we have

$$\begin{aligned} \int_0^u T'f(t) dt &= b \int \{e(u^{-1}t) - g(u^{-1}t)\}f(t) dt \\ &= b \int_0^u f(t) dt - b \int g(u^{-1}t)f(t) dt \\ &= \int_0^u \{bf(t) - V'f(t)\} dt \quad (\text{say}) \end{aligned}$$

for all  $u$ . Hence,  $T' = bI - V'$ . Setting  $b = \bar{a}$ , this implies that  $T = aI - V$ . Clearly,  $0 \in \sigma_c(V)$ . The other way implication is trivially true.

As an immediate consequence of the theorem we have:

**COROLLARY 4.** *Let  $V \in G_0$  be such that  $0 \in \sigma_p(V)$ . If  $T$  is invertible, and if  $T = aI - V$  for some scalar  $a$ , then  $\sigma(T) = \sigma_p(T)$ .*

**THEOREM 7.** *Let  $T$  be 0-adjoint, and let  $V = aI - T$  for some non-zero real  $a$ . If  $\sigma_c(V) \neq \emptyset$ , then  $\sigma(T) = \sigma_c(T)$ .*

**PROOF.** Clearly,  $V$  is 0-adjoint (with  $V' = V$ ). By Corollary 1, both  $\sigma_p(V)$  and  $\sigma_p(T)$  are symmetric about the origin. Assume that  $b \in \sigma_p(V)$ . Then there exist non-trivial functions  $f$  and  $g$  in  $L^p$  such that  $Vf = bf$  and  $Vg = -bg$ , or equivalently that  $(aI - V)f = (a - b)f$  and  $(aI - V)g = (a + b)g$ . This implies that  $(a + b) \in \sigma_p(T)$  and so also that  $\pm(a \pm b) \in \sigma_p(T)$ . But if  $-(a \pm b) \in \sigma_p(T)$ , then  $(2a + b) \in \sigma_p(V)$ . Repeating this argument we see that if  $b \in \sigma_p(V)$ , then  $(na + b) \in \sigma_p(V)$  for all even integers  $n$ . But this in view of the fact that  $\sigma(V)$  is a bounded subset of  $\mathbb{C}$  is absurd. Hence  $\sigma_p(V)$ , and so also  $\sigma_p(T)$  is empty. This completes the proof.

It is an immediate consequence of Theorems 6 and 7 that if the 0-adjoint  $T$  is invertible, then the spectrum of  $T$  consists either purely of the point spectrum of  $T$ , or purely of the continuous spectrum of  $T$ . Assuming now that  $\sigma(T) = \sigma_p(T)$  for the 0-adjoint invertible mapping  $T$ , we ask the question: Does  $\sigma_p(T)$  consist of a two point set?

**THEOREM 8.** *If  $i\beta \notin \sigma_p(TH - \beta H)$  for some  $\beta \in \varrho(T)$ , then  $0 \in \sigma_c(T)$ .*

**PROOF.** Suppose that  $i\beta \notin \sigma_p(TH - \beta H)$  for some  $\beta \in \varrho(T)$ . Then, since  $\sigma(TH - \beta H) = \sigma_p(TH - \beta H) \cup \sigma_c(TH - \beta H)$ , either  $i\beta \in \sigma_c(TH - \beta H)$  or  $i\beta \in \varrho(TH - \beta H)$ . Suppose that  $i\beta \in \varrho(TH - \beta H)$ . Then to each non-trivial  $f \in L^p$  there corresponds a unique non-trivial  $g \in L^p$  such that

$$(12) \quad (TH - \beta H)f - i\beta f = g.$$

Let  $f = f_1 + f_2$  and  $g = g_1 + g_2$ , where  $f_1, g_1 \in M_p$  and  $f_2, g_2 \in N_p$ . (The subspace  $M_p$  and  $N_p$  are as defined in Section 4.) We have from (12) that

$$(13) \quad i(T - \beta I)(f_1 - f_2) - i\beta(f_1 - f_2) = g_1 + g_2,$$

or

$$T(f_1 - f_2) - 2\beta f_1 = -i(g_1 - g_2).$$

Applying  $H$  to both sides of (12), we also have

$$(14) \quad T(f_1 + f_2) - 2\beta f_2 = i(g_1 - g_2).$$

From (13) and (14), we obtain

$$(15) \quad (T - \beta I)f_1 = (\beta f_2 - ig_2),$$

and so also (upon applying  $H$  to both sides of (15))

$$(16) \quad (T - \beta I)f_1 = -(\beta f_2 - ig_2).$$



Taken together, (15) and (16) imply that

$$(17) \quad (T - \beta I)f_1 = 0,$$

and so, since  $(T - \beta I)$  is invertible, that  $f_1 = 0$ . Letting  $f_1 = 0$  in (13) and (16) we have

$$(T - \beta I)f_2 = 0,$$

and hence that  $f_2 = 0$ . But then  $f = f_1 + f_2 = 0$ : this contradiction implies that  $i\beta \notin \rho(TH - \beta H)$ .

Now let  $i\beta \in \sigma_c(TH - \beta H)$ . Then to each  $f \in L^p$  there corresponds just one  $g \in L^p$  such that (12), and so also relations (13) to (17) hold. But then, by (17) and (13),

$$-Tf_1 - Tf_2 = -i(g_1 + g_2), \quad \text{or} \quad Tf = ig.$$

This implies that  $0 \in \sigma_c(T)$ .

## 6. Examples

In addition to the examples already considered during the course of Sections 4 and 5, in this section we consider certain further examples illustrating the preceding theory.

(1) Let the scalar  $b$  be as in Theorem 4, and let  $c$  and  $d$  be some scalars such that  $(c/d) \neq \pm \sqrt{b}$ . Let  $T^2 = bI$ , and let  $P = cI + dT$ . Then  $\sigma(P) = \sigma_p(P) = \{c \pm d\sqrt{b}\}$ . In particular, letting  $P = P_{(\alpha)} = \cos \pi\alpha + \sin \pi\alpha H$ ,  $0 < \alpha < 1$ , we see that  $\sigma(P_{(\alpha)}) = \sigma_p(P_{(\alpha)})$  lies on the unit circle. Mappings of the type  $P_{(\alpha)}$  have been considered by a number of authors, amongst them Juberg [7], Kober [8], Samko [10], and Duggal [6].

(2) Let  $S$  be the Stieltjes transform, and for  $b > 1$ , set  $T = bI - S$ . Then  $T$  satisfies the hypotheses of Theorem 7, and so  $\sigma(S) = \sigma_c(S)$ . Indeed,  $\sigma_c(S) = [0, 1]$  (see [11] and [12]).

(3) The mapping  $H$  is not essential to the validity of (a version of) Theorem 8: indeed, any  $T$  which satisfies  $T^2 = bI$  on  $L^p$  would do. For then  $\sigma(T) = \sigma_p(T) = \{\pm \sqrt{b}\}$ , and correspondingly we have a direct sum decomposition  $L^p = M_p \oplus N_p$ , where  $Tf = \sqrt{b}f$  for each  $f \in M_p$  and  $Tf = -\sqrt{b}f$  for each  $f \in N_p$ . We now employ this remark to give an example illustrating Theorem 8.

Let  $M_{(0)}$  (i.e.  $M_{(v)}$  with  $v=0$ ) be the (Meijer) transform defined in Section 3. Then  $T = M_{(0)}^2 \in G_0 \cap C(L^2(0, \infty))$ . Let  $P \in G_0$  be some invertible mapping such that  $P^2 = I$ . Then, denoting the Mellin transform of the function  $f \in L^2(0, \infty)$  by  $\hat{f}$ , we have that there exists a function  $K \in L^\infty$  such that

$$(Pf)^\wedge(y) = K(y)f^\wedge(y), \quad (K(y))^2 = 1, \quad y \in \mathbb{R}$$

(see [11]). Also, it is seen ([12]) that

$$(Tf)^\wedge(y) = (\pi/(\cosh \pi y + 1))f^\wedge(y).$$



The function  $\left(\pi - \frac{\pi}{\cosh \pi y + 1}\right)$  is boundedly invertible, and so it follows that  $\pi \in \rho(T)$ . Since  $(TPf - \pi Pf)^\wedge(y) = \pi f^\wedge(y)$ , or

$$\left(\frac{\pi}{\cosh \pi y + 1} - \pi\right) K(y) f^\wedge(y) = \pi f^\wedge(y),$$

i.e.

$$\left(\frac{-\cosh \pi y}{\cosh \pi y + 1}\right) f^\wedge(y) = K(y) f^\wedge(y)$$

holds a.e. only if  $f^\wedge(y) = 0$ , we see from Theorem 8 that  $0 \in \sigma_c(T)$ . Thus  $0 \in \sigma_c(M_{(0)})$ . The mapping  $T$  being 0-adjoint,  $\sigma(T) = \sigma_c(T)$ , and so  $\sigma(M_{(0)}) = \sigma_c(M_{(0)})$ . Indeed, it can be shown that  $\sigma(T) = \sigma_c(T) = [0, \pi/2]$ .

#### REFERENCES

- [1] CARTON-LEBRUN, C., Product properties of Hilbert transform, *J. Approximation Theory* **12** (1977), 356—360. *MR* **56** # 9173.
- [2] DUGGAL, B. P., Functional equations and linear transformations IIIA: Permutability and inversion, *Period. Math. Hungar.* **9** (1978), 93—107. *MR* **57** # 17059.
- [3] DUGGAL, B. P., Functiona equations and linear transformations IV: Interpolation, *Math. Ann.* **237** (1978), 277—285. *MR* **80a**: 39009.
- [4] DUGGAL, B. P., On the spectrum of a class of integral transforms, *J. Math. Anal. Appl.* **78** (1980), 41—48. (Corrigendum and Addendum, *ibid.* **95** (1983), p. 598.) *MR* **82b**: 44003.
- [5] DUGGAL, B. P., Spectrum of a class integral transforms, *Indian J. Pure Appl. Math.* **12** (1981), 964—970. *MR* **82i**: 44009.
- [6] DUGGAL, B. P., On a functional relation satisfied by fractional integrals, *Bull. London Math. Soc.* **15** (1983), 329—335.
- [7] JUBERG, R. K., The spectra for operators of a basic collection, *Bull. Amer. Math. Soc.* **79** (1973), 821—824. *MR* **48** # 4840.
- [8] KOBER, H., A modification of Hilbert transforms, the Weyl integral and functional equations, *J. London Math. Soc.* **42** (1967), 42—50. *MR* **34** # 3236.
- [9] POLLARD, H., Integral transforms, *Duke Math. J.* **13** (1946), 307—330. *MR* **8**—265.
- [10] SAMKO, S. G., Abel's generalised equation, Fourier transforms and convolution type equations, *Soviet Math. Dokl.* **12** (1971), 125—128. *MR* **40** # 4713.
- [11] DE SNOO, H. S. V., On the spectrum of Watson transforms, *J. London Math. Soc.* **8** (1974), 297—305. *MR* **50** # 1044.
- [12] TITCHMARSH, E. C., *Introduction to the Theory of Fourier Integrals*, Clarendon Press, Oxford, 1948.

(Received April 27, 1984)

SCHOOL OF MATHEMATICAL SCIENCES  
UNIVERSITY OF KHARTOUM  
P.O. BOX 321  
KHARTOUM  
SUDAN

# ON THE SPECTRUM OF A CLASS OF INTEGRAL TRANSFORMS III: AN APPLICATION

B. P. DUGGAL

## Abstract

Denote by  $t(a)$  the operator of dilatation by amount  $a$ ;  $a \in (-\infty, \infty)$ ,  $a \neq 0$ . We say that the continuous linear mapping  $T$  on  $L^p$  into itself belongs to the class  $G_0$  if it satisfies  $Tr(a) = m(a)t(a)T$ , where  $m(a) = (\text{sgn } a)$  or 1. Let  $1 < p, q < \infty$ . Using an explicit determination of the spectrum of mappings  $T \in G_0$ , it is shown that if  $T^2 = bI$  on  $L^p$ , then, for each  $f \in L^p$  and  $g \in L^q$ ,  $(1/p) + (1/q) < 1$ ,  $T$  satisfies the product relations  $T\{T(f)T(g) + bfg\} = b\{fT(g) + gT(f)\}$  and  $T\{fT(g) + gT(f)\} = T(f)T(g) + bfg$ . This generalizes a result which has been known for the Hilbert transform  $H$  for some time.  $H$  also satisfies the relation  $\sum_{r=1}^n H^r = 0$ . Let  $P = \alpha I + \beta T$  for some scalars  $\alpha, \beta$  ( $\beta \neq 0$ ).

We ask: Suppose that  $P$  satisfies  $\sum_{r=1}^n P^r = 0$  for some suitable integer  $n$ ; then does  $P$  or some power of  $P$  satisfy product relations of the type above? Again, if  $P$  satisfies product relations of the type above, then does  $P$  satisfy  $\sum_{r=1}^n P^r = 0$  for some suitable  $\alpha, \beta$  and  $n$ ? Let  $n = 4m$ ;  $m \geq 1$  some integer.

It is shown that if  $P$  is not scalar type, if 0 does not belong to the point spectrum of  $\sum_{r=1}^n P^r$  for all  $s < 4m$ , and if  $\sum_{r=1}^{4m} P^r = 0$ , then  $P^m$  satisfies product relations of the type above. Also, it is shown that if  $Q = i\sqrt{b}P$  satisfies the product relations above, and if neither of the points  $\pm\sqrt{b}$  is in the continuous spectrum of  $Q$ , then  $P$  satisfies  $\sum_{r=1}^{4m} P^r = 0$ .

## 1. Introduction

Let  $C(L^p)$ ,  $1 < p < \infty$ , denote the class of linear transformations which are continuous on  $L^p$  ( $= L^p(-\infty, \infty)$ ) into itself. The mapping  $T \in C(L^p)$  is said to belong to the class  $G_0$  if  $T$  satisfies the functional equation

$$(1) \quad Tt(a) = m(a)t(a)T, \quad -\infty < a < \infty, \quad a \neq 0,$$

where  $m(a) = (\text{sgn } a)$  or 1, and where  $t(a)$  is the operator of dilatation by amount  $a$ . It can be shown that if  $T$  satisfies functional equation (1), then there exists a Lebesgue measurable function  $k$  on  $(-\infty, \infty)$  such that

$$(2) \quad \int_0^u T f(t) dt = \int_{-\infty}^{\infty} m(x)k(ux^{-1})f(x) dx$$

(see [3], Lemma 2).

Let  $H$  denote the Hilbert transform

$$Hf(u) = (1/\pi)(P.V.) \int_{-\infty}^{\infty} (u-x)^{-1}f(x) dx.$$

1980 Mathematics Subject Classification. Primary 44A05; Secondary 47A10.

Key words and phrases.  $L^p$ -space, class  $G_0$  of integral transforms, spectrum, Hilbert transform.

Then  $H \in G_0$ . Cossar [2] and Tricomi [15] have shown that if  $f \in L^p$  and  $g \in L^q$ ,  $1 < p, q < \infty$  and  $(1/r) = (1/p) + (1/q) < 1$ , then  $H$  satisfies the product relations

$$(3) \quad \begin{aligned} H\{H(f)H(g) + bfg\} &= b\{fH(g) + gH(f)\}; \\ H\{fH(g) + gH(f)\} &= H(f)H(g) + bfg, \end{aligned}$$

where  $b = -1$ . Recently, Carton-Lebrun [1] has proved that the relations (3) remain valid even in the case in which  $r = 1$ . The proof given by Carton-Lebrun, which is different from that of [2] and [15] (see also [12]; p. 432), depends upon a judicious use of the Fourier transforms and the fact that  $H$  is weak type (1,1). We show here (see Section 3) that relations of the type (3) hold for continuously invertible elements  $T$  of the class  $G_0 \cap C(L^p)$  which satisfy the property that  $T^2 = bI$ , on  $L^p$ , for some non-zero scalar  $b$ .

Recall that  $HHf = -f$  for each  $f \in L^p$ . Hence  $H$  also satisfies the relation

$$(4) \quad \sum_{r=1}^n H^r = 0, \quad (n = 4),$$

on  $L^p$ . Let  $T \in G_0 \cap C(L^p)$ , and set, for some scalars  $\alpha, \beta$  ( $\beta \neq 0$ ),  $P = \alpha I + \beta T$  ( $I$  = the identity map). The question that we ask is the following. Suppose that, for some suitable choice of  $\alpha, \beta$  and  $n$ ,  $P$  satisfies (4). Then does  $P$  or some power of  $P$  satisfy relations of the type (3)? Again, suppose that  $P$  satisfies (3). Then does  $P$  satisfy (4) for some suitable  $n$ ? Let  $n = 4m$ ;  $m \geq 1$  some integer. We show (see Section 5) that if  $P^m$  is not scalar type (i.e.  $P^m \neq aI$  for some scalar  $a$ ), if 0 does not belong to the point spectrum of  $\sum_{r=1}^s P^r$  for all  $s < 4m$ , and if  $P$  satisfies  $\sum_{r=1}^{4m} P^r = 0$ , then  $P^m$  satisfies relations of the type (3). Also, we show that if  $Q = i\sqrt{b}P^m$  satisfies (3), and if neither of the points  $\pm\sqrt{b}$  is in the continuous spectrum of  $Q$ , then  $P^m$  satisfies (4).

We remark here that the Fourier transform technique used by Carton-Lebrun (or its counterpart — the Mellin transform technique — on the multiplicative group  $(0, \infty)$ ) seemingly cannot be extended to prove the analogue of (3) for mappings  $T \in G_0 \cap C(L^p)$ . Our technique below exploits the fact that an explicit determination of the spectrum  $\sigma(T)$  of mappings  $T$  under consideration can be made, and that  $\sigma(T) = \sigma_p(T)$  (= the point spectrum of  $T$ ) consists of a two point set.

## 2. Some notation

In addition to the notation already introduced, the following notation will be used in the sequel.  $\mathbf{R}$  will denote the set of reals, and  $\mathbf{C}$  will denote the set of complex numbers.  $\sigma_c(T)$  will denote the continuous spectrum of the mapping  $T$ . It is known (see [5]) that if  $T \in G_0 \cap C(L^p)$ , then  $\sigma(T) = \sigma_p(T) \cup \sigma_c(T)$ .  $L_p$  will denote  $L^p(0, \infty)$ , and  $p'$  will denote the index conjugate to  $p$  (i.e.  $(1/p) + (1/p') = 1$ ). Whenever the integration extends over all of  $\mathbf{R}$ , the limits will be omitted from the integrals. Although we do not always say so, all our equalities involving functions are to be considered as holding a.e. only. Any other notation will be defined as and when required.



## 3. Product relations

Let  $P = \alpha I + \beta T$ , where  $T \in G_0 \cap C(L^p)$ , and where  $\alpha, \beta$  ( $\beta \neq 0$ ) are some scalars.

THEOREM 1. If  $P^2 = bI$  on  $L^p$  for some non-zero scalar  $b$ , then

$$(I) \quad P\{P(f)P(g) + bfg\} = b\{fP(g) + gP(f)\};$$

$$(II) \quad P\{fP(g) + gP(f)\} = P(f)P(g) + bfg$$

for each  $f \in L^p$  and  $g \in L^q$ , where  $1 < p, q < \infty$  and  $(1/r) = (1/p) + (1/q) < 1$ .

PROOF. Since (I) and (II) hold trivially if  $T$  is scalar type, we assume in the following that contrary is the case. Then, since  $P^2 = bI$  on  $L^p$ , we have from [6, Theorem 4] that  $\sigma(P) = \sigma_p(P) = \{\pm\sqrt{b}\}$ . Let  $E_1$  be the projection associated with the spectral set  $\{+\sqrt{b}\}$  of  $P$ , and let  $E_2$  be the projection associated with the spectral set  $\{-\sqrt{b}\}$  of  $P$ . Let  $N_p$  and  $M_p$  ( $p$  as in  $L^p$ ) denote, respectively, the ranges of  $E_1$  and  $E_2$ . Then  $L^p = N_p \oplus M_p$ ; also,  $Pf = \sqrt{b}f$  for each  $f \in N_p$  and  $Pg = -\sqrt{b}g$  for each  $g \in M_p$ .

Let  $f \in L^p$  and  $g \in L^q$ . Then there exist unique functions  $f_1 \in N_p$  and  $f_2 \in M_p$ , and  $g_1 \in N_q$  and  $g_2 \in M_q$ , such that  $f = f_1 + f_2$  and  $g = g_1 + g_2$ . We have

$$\begin{aligned} P\{P(f)P(g) + bfg\} &= P\{P(f_1 + f_2)P(g_1 + g_2) + b(f_1 + f_2)(g_1 + g_2)\} = \\ &= bP\{(f_1 - f_2)(g_1 - g_2) + (f_1 + f_2)(g_1 + g_2)\} = \\ &= 2bP\{f_1 g_1 + f_2 g_2\}. \end{aligned}$$

Now,  $f_1 g_1 \in N_r$  and  $f_2 g_2 \in M_r$ . Hence

$$\begin{aligned} P\{P(f)P(g) + bfg\} &= 2b\{P(f_1 g_1) + P(f_2 g_2)\} = \\ &= 2b^{3/2}\{f_1 g_1 - f_2 g_2\} = \\ (5) \quad &= b\{(f_1 + f_2)(\sqrt{b}g_1 - \sqrt{b}g_2) + (g_1 + g_2)(\sqrt{b}f_1 - \sqrt{b}f_2)\} = \\ &= b\{fP(g) + gP(f)\}, \end{aligned}$$

i.e. (I) is satisfied.

It is clear that the resultant function on either side of (I) is in  $L^r$ . Hence, since  $P^2 = bI$  on  $L^r$  ( $1 < r < \infty$ ), (II) follows from (I) upon applying  $P$  to both sides.

The case  $r=1$ . The preceding argument does not extend to the case in which  $r=1$ , i.e.  $q=p'$ , unless some additional hypotheses are made on  $P$ . The difficulty here lies with relation (5). A set of conditions which ensure the validity of (5), and so also of (I), can be given as follows.

Let  $N_1$  ( $M_1$ ) denote the set of functions  $f$  (resp.,  $g$ ) such that  $f = f_1 f_2$  (resp.,  $g = g_1 g_2$ ) for some  $f_1 \in N_p$  and  $f_2 \in N_{p'}$  (resp.,  $g_1 \in M_p$  and  $g_2 \in M_{p'}$ ). The set  $A = N_1 \oplus M_1$  is then a subset of  $L^1$ . Assume that  $Pf = \sqrt{b}f$  for each  $f \in N_1$ , and that  $Pg = -\sqrt{b}g$  for each  $g \in M_1$  ( $b$  as in Theorem 1). Then for  $F \in A$  ( $\subset L^1$ ),  $TF \in A$  and relation (5), and so also (I), is satisfied. That relation (II) also holds follows from an argument similar to that used to prove (I). (Here one obviously cannot apply  $P$  to both sides of (I) to obtain (II).)

We prove now a converse to Theorem 1. An improved version of this converse is contained in Theorem 4 *infra*.

**THEOREM 2.** *Let  $T \in G_0 \cap C(L^p)$  be such that  $\sigma(T) = \sigma_p(T)$ . Define  $P = \alpha I + \beta T$ , where  $\alpha, \beta$  ( $\beta \neq 0$ ) are some scalars. If, for each  $f \in L^p$  and  $g \in L^q$  ( $1 < p, q < \infty$ ,  $(1/p) + (1/q) = (1/r) < 1$ ),  $P$  simultaneously satisfies (I) and (II) of Theorem 1, then  $P$  is invertible and  $P^2 = bI$  on  $L^p$ .*

**PROOF.** Since the theorem is trivially true in the case in which  $P$  is scalar type, we assume henceforth that contrary is the case. Clearly,  $\sigma(P) = \sigma_p(P)$ . From (I) and (II) we have that

$$\begin{aligned} P^2\{P(f)P(g) + bfg\} &= bP\{fP(g) + gP(f)\} = \\ &= b\{P(f)P(g) + bfg\}, \end{aligned}$$

or

$$(6) \quad (P^2 - b)\{P(f)P(g) + bfg\} = 0.$$

Since  $P(f)P(g) + bfg \in L^r$ , we have that  $b \in \sigma_p(P^2)$ . We now show that  $0 \notin \sigma_p(P)$ . For suppose that  $0 \in \sigma_p(P)$ . Then there exists a non-trivial  $f \in L^p$  such that  $Pf = 0$ . Fix  $f$ . Then we have from (I) that  $P(fg) = fP(g)$  for each  $g \in L^q$ , and so, since  $Pg \in L^q$ ,  $P^2(fg) = fP^2(g)$ . But, by (6),  $P^2(fg) = bfg$ ; so  $f(P^2g - bg) = 0$  for each  $g \in L^q$ . Hence  $P^2 = bI$  on  $L^q$ , and as such  $P$  is invertible on  $L^q$  for all  $1 < q < \infty$ . The contradiction implies that  $0 \notin \sigma_p(P)$ .

Now let  $a \in \sigma(P) = \sigma_p(P)$ . Then, upon choosing  $f \in L^p$  and  $g \in L^q$  such that  $Pf = af$  and  $Pg = ag$ , we have from (I) and (II) that  $2a/(a^2 + b) = (a^2 + b)/(2ab)$ , i.e.  $a^2 = b$ . Thus  $a \in \sigma_p(P)$  if and only if  $a^2 = b$ . Hence  $P^2 = bI$  on  $L^p$ , as required.

**REMARK.** (1) Let  $P$  be as in the statement of the theorem above. If  $P$  satisfies (I) and (II) of Theorem 1, then  $\sigma_p(P)$ , and so also  $\sigma_p(T)$ , cannot be empty. To see this, we note from (6) that  $b \in \sigma_p(P^2)$ . Hence either  $\sqrt{b}$  or  $-\sqrt{b}$  (or both)  $\in \sigma(P)$ . Suppose that  $\sqrt{b} \in \sigma_c(P)$ . (Recall that by Theorem 1 of [5],  $\sigma(P) = \sigma_p(P) \cup \sigma_c(P)$ .) Then to each  $F \in L^r$ ,  $(1/p) + (1/q) = (1/r) < 1$ , there corresponds just one  $G \in L^r$  such that  $(P - \sqrt{b})F = G$ , and so that  $(P^2 - b)F = \sqrt{b}G + PG$ . Choosing  $F = P(f)P(g) + bfg$  for some  $f \in L^p$  and  $g \in L^q$ , we have that  $PG = -\sqrt{b}G$ , i.e.  $-\sqrt{b} \in \sigma_p(P)$ . As a consequence of this observation we have that if  $P$  is such that  $\sigma(P) = \sigma_c(P)$ , then  $P$  cannot satisfy (I) and (II) of Theorem 1.

#### 4. Some consequences of Theorem 1

As already mentioned, the Hilbert transform  $H \in G_0 \cap C(L^p)$  and satisfies the identity  $HH = -I$  on  $L^p$  ( $1 < p < \infty$ ). It follows from Theorem 1 that relations (3) are satisfied (with  $b = -1$ ) for all  $(1/r) = (1/p) + (1/q) < 1$ . Since  $\sigma(H) = \sigma_p(H) = \{\pm i\}$ ,  $\{i\}$  and  $\{-i\}$  are spectral sets for  $H$ . The associated projections  $E_1$  and  $E_2$ , and their ranges, are identified as follows.

For each  $f \in L^p$ ,  $1 < p < \infty$ , define the mapping  $S$  by

$$Sf(z) = (1/(2\pi i)) \int (t-z)^{-1} f(t) dt, \quad z \in \mathbb{C} \setminus \mathbb{R}.$$

Then

$$(7) \quad S^+ f(x) = \lim_{y \rightarrow 0^+} Sf(x+iy), \quad S^- f(x) = \lim_{y \rightarrow 0^-} Sf(x+iy)$$

exist a.e., and satisfy

$$(8) \quad S^+ = (1/2i)(-H+iI), \quad S^- = (-1/2i)(H+iI), \quad S^+ - S^- = I$$

on  $L^p$  (see [1]). Identifying  $E_1$  with  $-S^-$  and  $E_2$  with  $S^+$ , we see that  $N_p = \{f \in L^p : Hf = if\}$  and  $M_p = \{f \in L^p : Hf = -if\}$ .

The set of functions  $f \in L^1$  which are representable by  $f = f_1 + f_2$ ,  $f_1 \in N_1$  and  $f_2 \in M_1$ , form a non-closed linear subspace  $A$  of  $L^1$ .  $A$  is identical with the set of functions  $f$  such that both  $f$  and  $Hf \in L^1$  (see [8], [9], [10]). It thus follows that:

**COROLLARY 1** (Carton-Lebrun [1]). *Relations (3) hold, with  $b = -1$ , for each  $f \in L^p$  and  $g \in L^q$ , where  $1 < p, q < \infty$  and  $(1/p) + (1/q) \leq 1$ .*

Let  $Q_\mu$  be the mapping  $Q_\mu f(x) = |x|^{-\mu} f(x)$ . The extended Hilbert transform  $H^{(\mu)}$ ,

$$H^{(\mu)} f(x) = Q_\mu (H Q_{-\mu} f)(x) = (1/\pi) |x|^{-\mu} (\text{P.V.}) \int \frac{|t|^\mu}{(x-t)} f(t) dt,$$

$\in G_0 \cap C(L^p)$  for all  $(-1/p') < \mu < (1/p)$ . Also  $H^{(\mu)} H^{(\mu)} f = -f$  for each  $f \in L^p$ ;  $1 < p < \infty$  (see [4]). Hence it follows that  $H^{(\mu)}$  satisfies (I) and (II), with  $b = -1$ , for all  $(1/p) + (1/q) < 1$  and  $\max((-1/p'), (-1/q')) < \mu < \min((1/p), (1/q))$ . Now set  $f_1 = Q_{-\mu} f$  and  $g_1 = Q_{-\mu} g$ . Then:

**COROLLARY 2.**

$$H Q_\mu \{H(f_1)H(g_1) - f_1 g_1\} = -Q_\mu \{f_1 H(g_1) + g_1 H(f_1)\};$$

$$H Q_\mu \{f_1 H(g_1) + g_1 H(f_1)\} = Q_\mu \{H(f_1)H(g_1) - f_1 g_1\}$$

for each  $Q_\mu f_1 \in L^p$  and  $Q_\mu g_1 \in L^q$ , where  $1 < p, q < \infty$ ,  $(1/p) + (1/q) < 1$  and  $\max((-1/p'), (-1/q')) < \mu < \min((1/p), (1/q))$ .

As another consequence of Theorem 1, we prove the following relation between mappings  $P, H$  and the boundary values of analytic functions.

**COROLLARY 3.** *Let  $P$  ( $\neq cI$  for some scalar  $c$ ) be as in Theorem 1. Let  $f \in L^p$  and  $g \in L^q$ , where  $1 < p, q < \infty$  and  $(1/p) + (1/q) < 1$ , be such that  $Pf = \sqrt{b}f$  and  $Tg = -\sqrt{b}g$ . Then*

$$fg = H(f)H(g), \quad P(f)P(g) = -bH(f)H(g).$$

Also

$$S^+(f)S^+(g) = S^-(f)S^-(g) = 0,$$

where  $S^+$  and  $S^-$  are as defined by (7).

PROOF. By Corollary 1 of [3],  $H$  and  $P$  commute, and so

$$PS^+ = (1/2i)(-PH + iP) = (1/2)(-HPH + iHP) = iHPS^+$$

on  $L^p$ . Similarly,

$$PS^- = -iHPS^-$$

on  $L^p$ . We have

$$\begin{aligned} P(f)P(g) + bfg &= P(S^+f - S^-f)P(S^+g - S^-g) + bfg = \\ &= HP(S^+f - S^-f)HP(S^+g - S^-g) + bfg = \\ &= HP(f)HP(g) + bfg = -bH(f)H(g) + fg. \end{aligned}$$

By (I) this implies that

$$-bP\{H(f)H(g) - fg\} = b\{fP(g) + gP(f)\} = -b^{3/2}\{fg - fg\} = 0,$$

i.e.  $H(f)H(g) = fg$ . Since  $-bfg = P(f)P(g)$ , we also have that  $P(f)P(g) = -bH(f)H(g)$ .

Since, by (8),  $H = -i(S^+ + S^-)$ , we have from  $H(f)H(g) = fg$  that

$$-(S^+f + S^-f)(S^+g + S^-g) = (S^+f - S^-f)(S^+g - S^-g),$$

i.e.

$$S^+fS^+g + S^-fS^-g = 0.$$

Applying  $H$  to both sides, and using the fact that  $H(S^+fS^+g) = -i(S^+fS^+g)$  and  $H(S^-fS^-g) = i(S^-fS^-g)$  (see [1]), we also have that

$$-(S^+fS^+g) + (S^-fS^-g) = 0.$$

This completes the proof.

## 5. Relations $\sum_{r=1}^n P^r = 0$

We tackle now the questions raised in paragraph three of the introduction. We assume in the following that our mappings  $T$  are not scalar type.

**THEOREM 3.** (a) Let the invertible mapping  $T \in G_0 \cap C(L^p)$  be such that  $T^2 = bI$  on  $L^p$  for some non-zero scalar  $b$ . Then there exist scalars  $\alpha, \beta$  ( $\beta \neq 0$ ) and an integer  $m \geq 1$  such that

(i)  $P^m$  is not scalar type;

(ii)  $0 \notin \sigma_p(\sum_{r=1}^n P^r)$  for all  $n < 4m$ ;

(iii)  $\sum_{r=1}^{4m} P^r = 0$ ,

where  $P = \alpha I + \beta T$ .

(b) If  $P$  satisfies (i), (ii) and (iii) above, then  $P$  is invertible and  $\sigma(P^m) = \sigma_p(P^m) = \{\pm i\}$ .



PROOF. (a) Since  $T^2 = bI$  on  $L^p$ , we have from [6, Theorem 4] that  $\sigma(T) = \sigma_p(T) = \{\pm \sqrt{b}\}$ . Choose  $\alpha = \cos \theta$ ,  $\beta = (i \sin \theta) / \sqrt{b}$ ,  $\theta = (\pi/2m)$  for some integer  $m \geq 1$ . Then  $P = \alpha I + \beta T = \cos \theta I + (i \sin \theta / \sqrt{b}) T = \cos \theta + \sin \theta V$  (say).

Clearly,  $V^2 = (-T^2/b) = -I$ , and  $\sigma(V) = \sigma_p(V) = \{\pm i\}$ . The sets  $\{+i\}$  and  $\{-i\}$  are spectral sets for  $V$ , and so there exist associated projections  $E_1$  and  $E_2$  such that  $E_1 + E_2 = I$  and  $iE_1 - iE_2 = V$ . Let  $N$  and  $M$  denote, respectively, the ranges of  $E_1$  and  $E_2$ . Then  $L^p = N \oplus M$ ,  $Vf = if$  for each  $f \in N$  and  $Vg = -ig$  for each  $g \in M$ .

Let  $f \in L^p$ . Then there exist unique  $f_1 \in N$  and  $f_2 \in M$  such that  $f = f_1 + f_2$ . We have

$$P^r f = P^r (f_1 + f_2) = P^{r-1} (Pf_1 + Pf_2) = P^{r-1} (e^{i\theta} f_1 + e^{-i\theta} f_2) = e^{ir\theta} f_1 + e^{-ir\theta} f_2.$$

Hence

$$\begin{aligned} \sum_{r=1}^n P^r f &= \sum_{r=1}^n \{e^{ir\theta} f_1 + e^{-ir\theta} f_2\} = \\ &= (e^{i\theta}(1 - e^{in\theta}) / (1 - e^{i\theta})) f_1 + (e^{-i\theta}(1 - e^{-in\theta}) / (1 - e^{-i\theta})) f_2 \end{aligned}$$

for each  $f = f_1 + f_2 \in L^p$ . Choosing  $n = 4m$ , we thus have that  $\sum_{r=1}^{4m} P^r f = 0$  for each  $f \in L^p$ . Since  $P^m f = i(f_1 - f_2)$ ,  $P^m$  is not scalar type. Also,  $0 \notin \sigma_p(\sum_{r=1}^n P^r)$  for all  $n < 4m$ .

(b) By (iii)  $P(f + Pf + \dots + P^{4m-1}f) = 0$  for each  $f \in L^p$ . Since  $0 \notin \sigma_p(\sum_{r=1}^n P^r)$  for all  $n < 4m$ ,  $0 \notin \sigma_p(P)$ . Hence  $\sum_{r=0}^{4m-1} P^r = 0$  ( $P^0 = I$ ) on  $L^p$ . This when taken in conjunction with (iii) implies that  $P^{4m} = I$ , i.e.  $P^{4m}$ , and so also  $P$ , is invertible. We now show that  $\sigma_c(P) = \emptyset$ .

Let  $a \in \sigma_c(P)$ . (Necessarily,  $a \neq 0$ .) Then to each  $f \in L^p$  there corresponds just one  $g \in L^p$  (the set of such  $g$  being a dense proper subset of  $L^p$ ) such that  $Pf - af = g$ . We have

$$(P^n - a^n)f = (P^{n-1} + aP^{n-2} + \dots + a^{n-2}P + a^{n-1}I)g.$$

Letting  $n = 4m$ , we have that  $P^{4m} = I$ ,  $a^{4m} = 1$ , and so (for a dense proper subset of  $L^p$ )

$$P^{4m-1} + aP^{4m-2} + \dots + a^{4m-2}P + a^{4m-1}I = 0,$$

or equivalently that

$$F(P) = aP^{4m-1} + a^2P^{4m-2} + \dots + a^{4m-1}P = -I.$$

Since  $\sigma(F(P)) = F(\sigma(P))$ , we now have that

$$(4m-1)a^{4m} = (4m-1) = -1,$$

i.e.  $m=0$  in (iii) on a dense proper subset of  $L^p$ . This contradiction implies that  $\sigma_c(P) = \emptyset$ , and so that  $\sigma(P) = \sigma_p(P)$ .

To conclude the proof, we now notice that  $\sigma(P^{2m}) = \sigma_p(P^{2m}) \subseteq \{\pm 1\}$ . However,  $1 \notin \sigma_p(P^{2m})$ , as the following argument shows. Suppose that  $\sigma_p(P^{2m}) = \{\pm 1\}$ . Then each  $f \in L^p$  may be uniquely written in the form  $f = f_1 + f_2$ , where  $P^{2m}f_1 = f_1$  and  $P^{2m}f_2 = -f_2$ . By (iii), we then have that

$$0 = \sum_{r=1}^{4m} P^r f = P(f_1 + f_2) + \dots + P^{2m-1}(f_1 + f_2) + (f_1 - f_2) + \\ + P(f_1 - f_2) + \dots + P^{2m-1}(f_1 - f_2) + (f_1 + f_2) = 2 \sum_{r=0}^{2m-1} P^r f_1.$$

This, however, contradicts (ii). Hence  $\sigma_p(P^{2m})$  is a proper subset of  $\{\pm 1\}$ . Since (iii) is trivially satisfied when  $P^{2m}f = -f$  (for each  $f \in L^p$ ), we conclude that  $P^{2m} = -I$ , and hence that  $\sigma(P^m) = \sigma_p(P^m) \subseteq \{\pm i\}$ . In view of (i),  $\sigma_p(P^m) = \{\pm i\}$ .

REMARKS. (2) The invertibility of  $P$ , where  $P = \alpha I + \beta T$  satisfies (i)–(iii), is not enough to guarantee the invertibility of the mapping  $T$ . Thus let

$$Af(x) = \alpha f(x) + \beta (\operatorname{sgn} x) (1/\pi^2) |x|^{-\mu} (\text{P.V.}) \int |t|^\mu \frac{\log t^2 - \log x^2}{(t-x)} f(t) dt.$$

where  $\alpha = i$ ,  $\beta = -i$ , and  $(-1/p') < \mu < (1/p)$ . Then  $A^2 = -I$  on  $(L^p)$  (see [5]),  $\sum_{r=1}^4 A^r = 0$  and  $0 \notin \sigma_p(\sum_{r=1}^n A^r)$  for all  $n < 4$ . However,  $T$  is not invertible: indeed,  $\sigma(T) = \sigma_p(T) = \{0, \pm 2\}$  (see [5]). The same example shows that the invertibility of  $T$  is not necessary to the validity of (i), (ii) and (iii) of the theorem.

(3) Let  $P$ , where  $P$  satisfies (iii) of the theorem, be one-one. Then  $P$  is invertible, and an argument similar to that used in the proof of (b) above shows that  $\sigma_c(P) = \emptyset$ . In consequence we have that if  $\sigma_c(T) \neq \emptyset$ , then  $P = \alpha I + \beta T$  does not satisfy (iii) above for any values of  $\alpha$ ,  $\beta$  ( $\beta \neq 0$ ) and  $m$ .

The following theorem depicts the close relation between relations of the types (3) and (4).

THEOREM 4. Let  $T \in G_0 \cap C(L^p)$ , and let  $P = \alpha I + \beta T$  for some scalars  $\alpha$ ,  $\beta$  ( $\beta \neq 0$ ). Then  $P$  satisfies (i), (ii) and (iii) of Theorem 3 if and only if  $Q = i\sqrt{b} P^m$  ( $b \neq 0$ ) satisfies (I), (II) of Theorem 1 and neither of the points  $\pm\sqrt{b} \in \sigma_c(Q)$ .

PROOF. If  $P$  satisfies (i)–(iii), then  $\sigma(P^m) = \sigma_p(P^m) = \{\pm i\}$ , and so  $\sigma(Q) = \sigma_p(Q) = \{\pm\sqrt{b}\}$ . By Theorem 1,  $Q$  satisfies (I) and (II).

Conversely, if  $Q$  satisfies (I) and (II), then as already seen,  $0 \notin \sigma_p(Q)$ . We show that  $0 \notin \sigma_c(Q)$ . For suppose that  $0 \in \sigma_c(Q)$ . Then there exist sequences of unit vectors  $\{f_n\} \in L^p$  and  $\{g_n\} \in L^q$  ( $q \neq p'$ ,  $1 < p, q < \infty$ ) such that  $\|Qf_n\| \rightarrow 0$  and  $\|Qg_n\| \rightarrow 0$ . By (I),

$$\|Q\{Q(f_n)Q(g_n) + bf_ng_n\}\| = \|b\|f_nQ(g_n) + g_nQ(f_n)\| \leq \\ \leq \|b\| \{\|f_n\| \|Q(g_n)\| + \|g_n\| \|Q(f_n)\|\} \rightarrow 0.$$

Since, as seen in the proof of Theorem 2,

$$Q^2\{Q(f_n)Q(g_n)+bf_ng_n\}=b\{Q(f_n)Q(g_n)+bf_ng_n\},$$

we have that

$$\|Q(f_n)Q(g_n)+bf_ng_n\| \rightarrow 0.$$

Again, since

$$\|Q(f_n)Q(g_n)+bf_ng_n\| \cong |b| \|f_n\| \|g_n\| - \|Q(f_n)\| \|Q(g_n)\|,$$

we have that

$$|b| \leq \lim_{n \rightarrow \infty} \|Q(f_n)Q(g_n)+bf_ng_n\| = 0.$$

Since  $b \neq 0$ , we have a contradiction. Thus  $0 \notin \sigma_c(Q)$ .

By the argument of Remark (1), if  $c \in \sigma_c(Q)$ , then  $-c \in \sigma_p(Q)$ . Since  $-c \in \sigma_p(Q)$  if and only if  $c^2=b$  (see the proof of Theorem 2), and since  $\pm\sqrt{b} \notin \sigma_c(Q)$ , we see that  $\sigma_c(Q)$  is empty, and that  $\sigma(Q)=\sigma_p(Q)=\{\pm\sqrt{b}\}$ . Thus  $\sigma(P^m)=\sigma_p(P^m)=\{\pm i\}$ . A simple argument now shows that (ii) and (iii) of Theorem 3 are also satisfied.

REMARK. (4) Let  $B \in G_0$  be 0-adjoint, i.e. the (Banach-space) adjoint mapping  $B'=dB$  for some scalar  $d$  such that  $|d|=1$  (see [3]). Then, if  $B$  is not scalar type,  $\sigma_p(B)$  is symmetric about the origin in the complex plane (see [6]). It follows that with  $P$  defined by 0-adjoint  $B$  the hypothesis that  $\pm\sqrt{b} \notin \sigma_c(P)$  may be omitted from Theorem 4. Can this be done generally?

## 6. Examples

(a) Both the Hilbert transform  $H$  and the extended Hilbert transform  $H^{(\mu)}$  satisfy (i), (ii), (iii), (I) and (II) with  $b=-1$ ,  $\alpha=0$ ,  $\beta=1$  and  $m=1$ .

(b) The mapping  $A$  of Remark (1) satisfies (i), (ii), (iii), (I) and (II) with  $m=1$  and  $b=-1$ . The mapping  $P=I-iA$  does not, however, satisfy (i), (ii) and (iii), or (I) and (II).

(c) The mapping  $P=(W'_\mu)^{-1}W_\mu=\cos \pi\mu+\sin \pi\mu H$ , where  $W_\mu$  denotes the fractional integral

$$W_\mu f(x) = (1/\Gamma(\mu)) \int_x^\infty (t-x)^{\mu-1} f(t) dt, \quad 0 < \mu < 1/p,$$

and  $W'_\mu$  denotes the mapping adjoint to  $W_\mu$ , has been studied by a number of authors (see [7], [11], [13]).  $P$  satisfies the hypotheses of Theorem 3 (a), and it is seen that for each  $\mu$  of the form  $\mu=1/(2m)$ ,  $p < 2m$ ,  $P$  satisfies (i), (ii) and (iii), and that  $P^m$  satisfies (I) and (II) with  $b=-1$ .

(d) Let  $M \in C(L_2)$  be the Meijer transform

$$Mf(x) = (\sqrt{2/\pi}) \int_0^\infty \sqrt{(xt)} K_\nu(xt) f(t) dt, \quad |\operatorname{Re} \nu| < 1.$$

Then the mapping  $T = M^2 \in G_0 \cap C(L_2)$ . The mapping

$$P = I - \frac{2 \cos \pi v}{\pi} M^2, \quad \operatorname{Im} v = 0, \quad \operatorname{Re} v \neq 0, \quad v \neq \pm 1/2,$$

is invertible: indeed, denoting the Mellin transform of  $f \in L_2$  by  $\hat{f}$ ,

$$(Pf)^{\wedge}(y) = \frac{\cosh \pi y - \cos \pi v}{\cosh \pi y + \cos \pi v} \hat{f}(y), \quad y \in \mathbf{R},$$

so that  $\sigma(P) = \sigma_c(P) = [\tan^2(\pi v/2), 1]$  if  $-1/2 < v < 1/2$  and  $\sigma(P) = \sigma_c(P) = [1, \tan^2(\pi v/2)]$  if  $-1 < v < -1/2$  or  $1/2 < v < 1$  (see [14]). Since  $\sum_{r=1}^n P^r f = 0$  for each  $f \in L_2$  if and only if  $\sum_{r=1}^n (P^r f)^{\wedge}(y) = 0$ , i.e. if and only if

$$\sum_{r=1}^n \left( \frac{\cosh \pi y - \cos \pi v}{\cosh \pi y + \cos \pi v} \right)^r = 0, \quad y \in \mathbf{R},$$

we see that  $P$  does not satisfy (iii). By Remark (1),  $P$  cannot satisfy (I) and (II).

(e) The Stieltjes transform  $S$ ,

$$Sf(x) = (1/\pi) \int_0^{\infty} (x+t)^{-1} f(t) dt, \quad 0 < x < \infty,$$

$\in G_0 \cap C(L_p)$ . Since  $\sigma(S) = \sigma_c(S) = [0, 1]$ , we have from Remark (1) that  $S$  does not satisfy (I) and (II). Now let  $p=2$ , and let, as in Example (d) above,  $\hat{f}$  denote the Mellin transform of  $f \in L_2$ . Let  $P = \alpha I - \beta S$ , where  $(\alpha/\beta) \notin [0, 1]$ . Then  $P$  is invertible, and

$$\left( \sum_{r=1}^n P^r f \right)^{\wedge}(y) = \sum_{r=1}^n \left( \frac{\alpha \cosh y - \beta}{\cosh y} \right)^r \hat{f}(y), \quad y \in \mathbf{R}.$$

Since the sum on the right-hand side in the above equality does not equal to zero for any positive integer  $n$  (and all  $y \in \mathbf{R}$ ), we see that  $P$  does not satisfy (iii).

#### REFERENCES

- [1] CARTON-LEBRUN, C., Product properties of Hilbert transform, *J. Approximation Theory* **21** (1977), 356—360. *MR* **56** # 9173.
- [2] COSSAR, J., On conjugate functions, *Proc. London Math. Soc.* **45** (1939), 369—381. *MR* **1**—52.
- [3] DUGGAL, B. P., Functional equations and linear transformations IIIA: Permutability and inversion, *Period. Math. Hungar.* **9** (1978), 93—107. *MR* **57** # 1703.
- [4] DUGGAL, B. P., Functional equations and linear transformations IV: Interpolation, *Math. Ann.* **237** (1978), 277—285. *MR* **80a**: 39009.
- [5] DUGGAL, B. P., On the spectrum of a class of integral transforms, *J. Math. Anal. Appl.* **78** (1980), 41—48. *MR* **82i**: 44009.
- [6] DUGGAL, B. P., On the spectrum of a class of integral transforms II, *Studia Sci. Math. Hungar.* **20** (1985), 451—460.
- [7] JUBERG, R. K., On the boundedness of certain singular integral operators, *Colloq. Math. Soc. János Bolyai Vol. 5, Hilbert Space Operators and Operator Algebras*, North-Holland, 1972, 305—318. *MR* **50** # 14374.

- [8] KOBER, H., A note on Hilbert transforms, *Quart. J. Math. Oxford* **14** (1943), 49—54. *MR* **5**—179.
- [9] KOBER, H., A note on Hilbert's operator, *Bull. Amer. Math. Soc.* **48** (1942), 421—426. *MR* **4**—40.
- [10] KOBER, H., An operator related to Hilbert transforms and to Dirichlet's integral, *J. London Math. Soc.* **39** (1964), 649—656. *MR* **29** # 3836.
- [11] KOBER, H., A modification of the Hilbert transform, the Weyl integral and functional equations, *J. London Math. Soc.* **42** (1967), 42—50. *MR* **34** # 3236.
- [12] OKIKIOLU, G. O., *Aspects of the Theory of Bounded Integral Operators in  $L^p$ -Spaces*, Academic Press, London—New York, 1971. *MR* **56** # 3581.
- [13] SAMKO, S. G., Operators of potential type, *Soviet. Math. Dokl.* **12** (1971), 125—128. *MR* **42** # 8340.
- [14] DE SNOO, H. S. V., On the spectrum of Watson transforms, *J. London Math. Soc.* **8** (1974), 297—305. *MR* **50** # 1044.
- [15] TRICOMI, F. G., *Integral Equations*, Interscience, New-York, 1965.

(Received April 27, 1984)

SCHOOL OF MATHEMATICAL SCIENCES  
UNIVERSITY OF KHARTOUM  
P.O. BOX 321  
KHARTOUM  
SUDAN



# ON COMPACT PACKING OF CIRCLES

A. FLORIAN

*To my brother Helmut Florian on his 60th birthday*

We shall concern ourselves with packings of circles on the two-dimensional unit sphere,  $S^2$ , and in the hyperbolic plane,  $H^2$ . By a circle on  $S^2$  we mean a spherical cap of radius less than  $\pi/2$ . A collection of closed circles is said to form a *packing* if no two of them have an inner point in common. Two circles of a packing with a common boundary point are called *neighbours*. Let  $O$  and  $O_i$  denote the centres of the circles  $c$  and  $c_i$ , respectively. Following L. Fejes Tóth [3], we define a packing of circles to be *compact* if each circle  $c$  of the packing satisfies the following three conditions:

- (i)  $c$  has a finite number of neighbours,
- (ii) if  $c$  has  $n$  neighbours,  $c_1, \dots, c_n$ , they can be numbered so that  $c_1$  touches  $c_2, \dots, c_n$  touches  $c_1$ ,
- (iii)  $c$  is covered by the union of the triangles  $OO_1O_2, \dots, OO_nO_1$ .

In [1, 3] compact packings in the Euclidean plane are considered. In particular, it is proved that the lower density of a compact packing of circles with radii from a given interval  $[a, b]$ , where  $0 < a < b < \infty$ , is at least  $\pi/\sqrt{12}$ . The bound  $\pi/\sqrt{12}$  is attained by the system of the in-circles of a regular tiling (6, 3). In the present paper we shall establish two analogous results concerning packings on the sphere and in the hyperbolic plane.

We shall use the same letter to denote a set and its area. If  $\{c_0, c_1, \dots\}$  is a packing of circles and  $S$  a bounded closed subset of  $S^2$  or  $H^2$ , the ratio  $\sum_i (c_i \cap S)/S$  is called the *density* of  $\{c_0, c_1, \dots\}$  with respect to  $S$  and is simply called the density of  $\{c_0, c_1, \dots\}$  if  $S = S^2$ . To formulate our result we define a function  $D(R)$  for  $R > 0$  (and  $R \leq \pi/3$  in the spherical case). Let  $d(R)$  be the density of three mutually touching circles of radius  $R$  with respect to the triangle spanned by their centres.

**THEOREM 1.** *If  $d$  is the density of a compact packing of a finite number of circles on  $S^2$  with radii not greater than  $R^* \leq \pi/3$  then*

$$(1) \quad d \cong d(R^*).$$

*Equality occurs in (1) if and only if  $\cos R^* = \frac{1}{2} \operatorname{cosec} \frac{\pi}{p}$ , where  $p = 2, 3, 4, 5$ , and the circles are the in-circles of the regular tiling  $(p, 3)$ .<sup>1</sup>*

<sup>1</sup> Theorem 1 confirms a conjecture by L. Fejes Tóth who drew the author's attention to this subject.

1980 *Mathematics Subject Classification*. Primary 52A45; Secondary 52A40.

*Key words and phrases*. Packing of circles, density, regular tiling.



It is well-known that, in contrast to the Euclidean plane, there is no satisfactory way of defining the density of a packing of circles with respect to the whole hyperbolic plane [2]. Thus we are restricted to proving a result of a 'local' character.

**THEOREM 2.** *In  $H^2$  let  $\{c_0, c_1, \dots\}$  be a compact packing of circles of radii not less than  $R^* > 0$ . There is a tiling  $\mathcal{T}$  with triangular faces and with its vertices at the centres of the circles so that the density of  $\{c_0, c_1, \dots\}$  with respect to any face of  $\mathcal{T}$  is at least  $d(R^*)$ .*

Let  $\cosh R^* = \frac{1}{2} \operatorname{cosec} \frac{\pi}{p}$ , where  $p = 7, 8, \dots$ , and let  $\{c_0, c_1, \dots\}$  be the system of the in-circles of the regular tiling  $(p, 3)$ . The centres of the circles are the vertices of the dual tiling  $(3, p)$ . The density of  $\{c_0, c_1, \dots\}$  with respect to any face of  $(3, p)$  is equal to  $d(R^*)$ .

**PROOF OF THEOREM 1.** Let  $c_0$  be any circle of a compact packing, and let  $c_1, \dots, c_n$  be the successive neighbours of  $c_0$ . Again,  $O_i$  denotes the centre of  $c_i$ . If  $O_0, O_1$  and  $O_2$  lie on the same great circle,  $c_0, c_1$  and  $c_2$  have radius  $\pi/3$ . Thus  $R^* = \pi/3$ , and Theorem 1 is obvious in this case. Consequently, in the following we may suppose that  $O_0, O_1$  and  $O_2$  (more generally  $O_0, O_i$  and  $O_{i+1}$ ) are the vertices of a non-degenerate triangle  $\Delta$ .

Let  $T_i$  and  $T_{12}$  be the touching points of the pairs  $(c_0, c_i)$  and  $(c_1, c_2)$  for  $i = 1, \dots, n$ . Let  $T_1T_2$  be the arc on the boundary of  $c_0$  which is covered by  $\Delta$ , and define  $T_2T_{12}$ ,  $T_{12}T_1$  and  $T_iT_{i+1}$  similarly. The arc-sided triangle bounded by  $T_1T_2$ ,  $T_2T_{12}$  and  $T_{12}T_1$  will be denoted by  $D$ . Observe that no neighbour of  $c_0$  touches  $c_0$  at an inner point of  $T_1T_2$ . Otherwise all neighbours of  $c_0$  with the exception of  $c_1$  and  $c_2$  are contained in  $D$  and therefore also in  $\Delta$ . But this contradicts condition (iii) of the definition of a compact packing. Thus the union of the arcs  $T_1T_2, \dots, T_nT_1$  covers the boundary of  $c_0$  without overlapping.

We proceed to prove that the triangles of type  $\Delta$  form a tiling.

Let  $U$  be the union of those circles of the packing that are contained in  $D$ . We assert that  $U$  is empty. Assuming the contrary, we choose a point  $P \in U \cap \operatorname{int} D$  and join  $P$  with an inner point  $Q$  of the arc  $T_1T_2$  by a Jordan curve  $l$  not intersecting  $c_1$  or  $c_2$  (Fig. 1). Since  $U$  is closed and  $Q$  does not belong to  $U$  we may assume that

$$(2) \quad U \cap l = \{P\}.$$

$P$  is a boundary point of some circle  $c'$  of the packing. Let  $c'_1, \dots, c'_m$  be the successive neighbours of  $c'$ , and  $T'_1, \dots, T'_m$  the touching points of these circles with  $c'$ , and let  $P \in T'_1T'_2$ .  $l$  intersects the interior of the arc-sided triangle  $D' = T'_1T'_2T'_{12}$ , where  $T'_{12}$  is the common point of  $c'_1$  and  $c'_2$ . Since  $D' \subset D$  and  $Q \notin U \cup c_1 \cup c_2$ ,  $Q$  is outside  $D'$ , so that  $l$  intersects the boundary of  $D'$  at a point  $X \neq P$ . From  $l \cap (c_1 \cup c_2) = \emptyset$  it follows that  $X \in U$  which contradicts (2). This shows that in fact no circle of the packing is contained in  $D$ . Therefore  $c_0$  and  $c_1$  are successive neighbours of  $c_2$ , and  $c_0$  and  $c_2$  are successive neighbours of  $c_1$ . This implies that any two different triangles of type  $\Delta$  do not overlap.

It is easy to see that each point  $Y \in S^2$  belongs to some triangle  $\Delta$ . In view of condition (iii) we may assume that  $Y \notin \bigcup_i c_i$ . Let  $Z$  be a point of  $\bigcup_i c_i$  at minimum

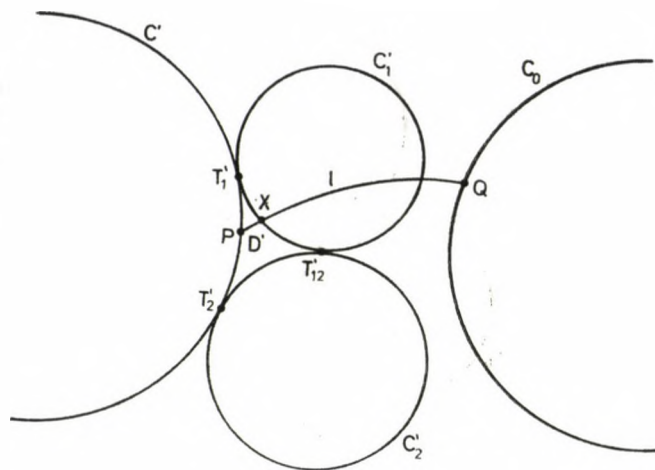


Fig. 1

distance from  $Y$ , say  $Z \in c_0$ . Then  $Y$  lies in one of the triangles  $O_0O_1O_2, \dots, O_0O_nO_1$ . Thus the triangles  $\Delta$  are the faces of a tiling  $\mathcal{T}$  as required.

Let  $\Delta = O_0O_1O_2$  be any face of  $\mathcal{T}$ , and let  $c_i$  be the circle of the packing centred at  $O_i$  ( $i=0, 1, 2$ ). Since  $c_3, c_4, \dots$  do not intersect  $\Delta$ , the density of  $\{c_0, c_1, c_2, \dots\}$  with respect to  $\Delta$  is

$$(3) \quad \frac{\sum_i (c_i \cap \Delta)}{\Delta} = \frac{\sum_{i=0}^2 (c_i \cap \Delta)}{\Delta}.$$

To prove Theorem 1 it suffices to show that

$$(4) \quad \frac{\sum_{i=0}^2 (c_i \cap \Delta)}{\Delta} \cong d(R^*)$$

for any three mutually touching circles  $c_0, c_1, c_2$  of radii not greater than  $R^* \leq \pi/3$ . Equality occurs in (4) if and only if the circles have radius  $R^*$ .

Let  $R_i$  be the radius of  $c_i$  ( $i=0, 1, 2$ ). Let  $T_1, T_2, T_{12}$  be the points of tangency of the pairs  $(c_0, c_1)$ ,  $(c_0, c_2)$  and  $(c_1, c_2)$  (Fig. 2). We denote the centre and the radius of the in-circle of  $\Delta$  by  $O$  and  $r$ . The in-circle touches the sides of  $\Delta$  at  $T_1, T_2, T_{12}$ . Writing  $2\beta_i$  ( $i=0, 1, 2$ ) for the angles of  $\Delta$ , and  $2\alpha_0, 2\alpha_1, 2\alpha_2$  for the angles  $T_2OT_1, T_1OT_{12}, T_{12}OT_2$  we have

$$(5) \quad \sum_{i=0}^2 \alpha_i = \pi, \quad \Delta = 2 \sum_{i=0}^2 \beta_i - \pi, \quad c_i \cap \Delta = 2\beta_i(1 - \cos R_i),$$

$$(6) \quad \cos \beta_i = \cos r \sin \alpha_i, \quad \tan R_i = \sin r \tan \alpha_i \quad (i = 0, 1, 2).$$

We consider an equilateral triangle  $\bar{\Delta}$  with in-radius  $r$ , and three congruent circles  $\bar{c}_0, \bar{c}_1, \bar{c}_2$  touching one another and having their centres at the vertices of  $\bar{\Delta}$ . If  $\bar{c}_i$



we have

$$f(\alpha) = f_1 \left( 1 - q + \frac{1}{\cos r} f_1' \right),$$

whence

$$(15) \quad f''(\alpha) = (1 - q)f_1'' + \frac{3}{\cos r} f_1' f_1'' + \frac{1}{\cos r} f_1 f_1'''.$$

From (14) it follows that

$$(16) \quad f_1''(\alpha) = \sin^2 r \cos r \frac{\sin \alpha}{(1 - \cos^2 r \sin^2 \alpha)^{3/2}}$$

and

$$(17) \quad f_1'''(\alpha) = \sin^2 r \cos r \frac{(1 + 2 \cos^2 r \sin^2 \alpha) \cos \alpha}{(1 - \cos^2 r \sin^2 \alpha)^{5/2}}.$$

By use of (11), (14), (16) and (17) we find

$$f_1'(\alpha) = -\frac{\cos r \cos \alpha}{\sin \beta}, \quad f_1''(\alpha) = \frac{\sin^2 r \cos \beta}{\sin^3 \beta},$$

$$f_1'''(\alpha) = \sin^2 r \cos r \frac{(1 + 2 \cos^2 \beta) \cos \alpha}{\sin^5 \beta},$$

so that, by (15),

$$(18) \quad \frac{\sin^4 \beta}{\sin^2 r \cos \beta} f''(\alpha) = (1 - q) \sin \beta + \cos \alpha \left( \beta \frac{1 + 2 \cos^2 \beta}{\sin \beta \cos \beta} - 3 \right).$$

Because of

$$\frac{1 + 2 \cos^2 \beta}{\sin \beta \cos \beta} = 2 \frac{2 + \cos 2\beta}{\sin 2\beta}$$

the second term on the right-hand side of (18) is proved to be positive for  $0 < \beta < \pi/2$  if

$$g(x) = x - 3 \frac{\sin x}{2 + \cos x}$$

is positive for  $0 < x < \pi$ . But this follows immediately from  $g(0) = 0$  and

$$g'(x) = \frac{(1 - \cos x)^2}{(2 + \cos x)^2} > 0.$$

By (7),  $1 - q$  is positive, too. Thus

$$(19) \quad f''(\alpha) > 0,$$

so that  $f(\alpha)$  is strictly convex for  $0 < \alpha < \pi/2$ . Applying Jensen's inequality we

obtain from (5) and (10)

$$\begin{aligned}
 \sum_{i=0}^2 (c_i \cap \Delta) &= 2 \sum_{i=0}^2 \beta_i (1 - q - \cos R_i) + 2q \sum_{i=0}^2 \beta_i = \\
 (20) \qquad \qquad &= 2 \sum_{i=0}^2 f(\alpha_i) + (\Delta + \pi)q \equiv 6f\left(\frac{\pi}{3}\right) + (\Delta + \pi)q.
 \end{aligned}$$

If  $2\bar{\beta}$  denotes the angle of  $\bar{\Delta}$ , then by (5) and (7)

$$f\left(\frac{\pi}{3}\right) = \bar{\beta}(1 - q - \cos \bar{R}), \quad q = \frac{6\bar{\beta}(1 - \cos \bar{R})}{6\bar{\beta} - \pi},$$

whence

$$6f\left(\frac{\pi}{3}\right) = -6\bar{\beta}q + (6\bar{\beta} - \pi)q = -\pi q.$$

Thus (20) implies (8). Equality holds only if  $\Delta = \bar{\Delta}$ .

Because of  $\pi/3 \leq \max \{\alpha_0, \alpha_1, \alpha_2\}$  we have by (6)

$$\bar{R} \leq \max \{R_0, R_1, R_2\} \leq R^*.$$

We establish (9) by showing that  $q = D(\bar{R})$  is a strictly decreasing function for  $0 < \bar{R} \leq \pi/3$ . (6) implies the relation

$$(21) \qquad \qquad \qquad \cos \bar{R} = \frac{1}{2 \sin \bar{\beta}},$$

so that

$$(22) \qquad \qquad q = F(\bar{\beta}) = \frac{6\bar{\beta} \left(1 - \frac{1}{2 \sin \bar{\beta}}\right)}{6\bar{\beta} - \pi}.$$

By (21) it suffices to prove that

$$F(x) = \frac{x \left(1 - \frac{1}{2 \sin x}\right)}{x - \frac{\pi}{6}}$$

is strictly decreasing for  $0 < x \leq \pi/2$ . This is an immediate consequence of the fact that

$$F_1(x) = 2x - \frac{x}{\sin x}$$

is strictly concave for  $0 < x \leq \pi/2$ . From

$$-2F_1''(x) = \frac{3x + x \cos 2x - 2 \sin 2x}{\sin^3 x}$$

we see that we have to show that

$$F_2(y) = y - \frac{4 \sin y}{3 + \cos y} > 0$$

for  $0 < y < \pi$ . This is obvious in view of  $F(0)=0$  and

$$F'_2(y) = \frac{(5 - \cos y)(1 - \cos y)}{(3 + \cos y)^2} > 0$$

and completes the proof of (4) and (1). Equality occurs in (1) if and only if the faces of the tiling  $\mathcal{T}$  are congruent equilateral triangles of side length  $2R^*$ . In this case  $\mathcal{T}$  is a regular tiling  $(3, p)$ , and  $R^*$  is equal to  $\arccos\left(\frac{1}{2} \operatorname{cosec} \frac{\pi}{p}\right)$ .

PROOF OF THEOREM 2. Let us construct a tiling  $\mathcal{T}$  with triangular faces and with its vertices at the centres of the circles, as in the proof of Theorem 1. A face  $\Delta$  of  $\mathcal{T}$  intersects only those circles of the packing that have their centres at the vertices of  $\Delta$ . If these circles are  $c_0, c_1, c_2$ , the density of  $\{c_0, c_1, \dots\}$  with respect to  $\Delta$  is given by (3). It remains to show that (4) applies to any three mutually touching circles of radii not less than  $R^* > 0$ .

The proof is quite analogous to that given in the spherical case. Keeping up the notations introduced above we have only to modify some formulae following (4). Equations (5) and (6) are to be replaced by

$$(5') \quad \sum_{i=0}^2 \alpha_i = \pi, \quad \Delta = \pi - 2 \sum_{i=0}^2 \beta_i, \quad c_i \cap \Delta = 2\beta_i (\cosh R_i - 1),$$

$$(6') \quad \cos \beta_i = \cosh r \sin \alpha_i, \quad \tanh R_i = \sinh r \tan \alpha_i \quad (i = 0, 1, 2).$$

To prove (8) and (9) we define a function  $f(\alpha)$  for  $0 < \alpha < \hat{\alpha} = \arcsin \frac{1}{\cosh r}$  by

$$(10') \quad f(\alpha) = \beta (\cosh R + q - 1),$$

where  $\beta$  and  $R$  are connected with  $\alpha$  by

$$(11') \quad \cos \beta = \cosh r \sin \alpha, \quad \tanh R = \sinh r \tan \alpha \quad (0 < \beta < \pi/2),$$

and  $q$  is given by (7). By using (11') we obtain from (10')

$$(18') \quad \frac{\sin^4 \beta}{\sinh^2 r \cos \beta} f''(\alpha) = (1 - q) \sin \beta + \cos \alpha \left( \beta \frac{1 + 2 \cos^2 \beta}{\sin \beta \cos \beta} - 3 \right).$$

Since the right-hand side of (18') was proved to be positive, we again have (19), so that  $f(\alpha)$  is strictly convex for  $0 < \alpha < \hat{\alpha}$ . Observing that, by (6'),  $\max \{\alpha_0, \alpha_1, \alpha_2\} < \hat{\alpha}$ , and applying Jensen's inequality, we obtain from (5') and (10')

$$(20') \quad \begin{aligned} \sum_{i=0}^2 (c_i \cap \Delta) &= 2 \sum_{i=0}^2 \beta_i (\cosh R_i + q - 1) - 2q \sum_{i=0}^2 \beta_i = \\ &= 2 \sum_{i=0}^2 f(\alpha_i) - (\pi - \Delta)q \geq 6f\left(\frac{\pi}{3}\right) - (\pi - \Delta)q. \end{aligned}$$

The required inequality (8) now follows from (20') and the relations

$$f\left(\frac{\pi}{3}\right) = \beta(\cosh \bar{R} + q - 1), \quad q = \frac{6\beta(\cosh \bar{R} - 1)}{\pi - 6\beta}.$$

Equality holds only if  $\Delta = \bar{\Delta}$ .

Because of  $\min\{\alpha_0, \alpha_1, \alpha_2\} \leq \pi/3$ , we have by (6')

$$R^* \leq \min\{R_0, R_1, R_2\} \leq \bar{R}.$$

We establish (9) by showing that  $q = D(\bar{R})$  is a strictly increasing function for  $\bar{R} > 0$ . Since by (6')

$$(21') \quad \cosh \bar{R} = \frac{1}{2 \sin \beta},$$

it remains to prove that

$$(22') \quad q = \frac{6\beta \left( \frac{1}{2 \sin \beta} - 1 \right)}{\pi - 6\beta}$$

is strictly decreasing for  $0 < \beta < \pi/6$ . The right-hand side of (22') is identical with the function  $F(\beta)$  defined by (22) which was proved to be strictly decreasing for  $0 < \beta \leq \pi/2$ .

(8) and (9) imply (4) with equality if and only if the circles have radius  $R^*$ . This concludes the proof of Theorem 2.

From (21), (22) and (21'), (22') we infer that  $\lim_{R^* \rightarrow 0} D(R^*) = \pi/\sqrt{12}$  which is the precise lower bound for the density of a compact packing of circles in the Euclidean plane.

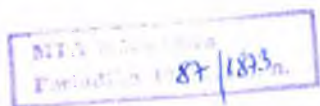
I would like to thank my son Andreas for drawing the figures.

#### REFERENCES

- [1] BEZDEK, A., BEZDEK, K. and BÖRÖCZKY, K., On compact packings, *Studia Sci. Math. Hungar.* **21** (1986).
- [2] BÖRÖCZKY, K., Sphere packings in spaces of constant curvature, *Mat. Lapok* **25** (1974), 265—306 (in Hungarian).
- [3] FEJES TÓTH, L., Compact packing of circles, *Studia Sci. Math. Hungar.* **19** (1984), 103—107.

(Received May 4, 1984)

INSTITUT FÜR MATHEMATIK  
UNIVERSITÄT SALZBURG  
HELLBRUNNERSTRASSE 34  
A-5020 SALZBURG  
AUSTRIA





## BOOK REVIEWS

**Skornyakov, L. A., Elements of General Algebra** (in Russian), Moscow, Nauka, 1983, 272 pp.

The aim of this book is to give an introduction to the most important chapters of abstract algebra. The author is a well-known algebraist, a leading professor of the Lomonosov State University (Moscow). We list the detailed contents of the book: Chapter I. Partially ordered sets and complete lattices 23 pp. Chapter II. Universal algebras (operations, algebras, congruences, varieties, free universal algebras) 35 pp. Chapter III. Lattices (modular and distributive lattices, Boolean algebras) 20 pp. Chapter IV. Associative rings and modules (von Neumann regular rings, Noetherian rings, tensor product, prime rings, radical and classical semisimple rings, Wedderburn—Artin Theorem, and homological algebra) 54 pp. Chapter V. Groups and Lie algebras (subgroups of free groups, nilpotent groups, linear groups, rings and Lie algebras, nilpotent Lie algebras and nilpotent groups) 38 pp. Chapter VI. Fields and skew fields (structure of field extensions, the fundamental theorem of Galois theory, Theorem of Frobenius concerning finite dimensional algebras over the real field) 30 pp. Chapter VII. Algebras with complemented lattice (in the sense of Bourbaki): ordered groups, normed rings (in the sense of Gelfand), topological rings 32 pp. Chapter VIII. Categories: fundamental concepts, additive categories 26 pp. There are interesting exercises and bibliography at the end of each chapter. There is a non-empty intersection with the author's earlier book entitled *Complemented Modular Lattices and Regular Rings* (in Russian), Moscow, Nauka, 1961. It would be a very good idea to publish the English translation of this modern monograph.

*S. Lajos*

**Vincze, I., Mathematische Statistik mit industriellen Anwendungen**, 2. ed., Bibliographisches Institut AG, Mannheim and Publishing House of the Hungarian Academy, Budapest, 1984, 502 pp. ISBN 963 05 3351 0.

This is the enlarged and improved second edition of the German version of the book; it is a joint edition of the Bibliographisches Institute AG and the Publishing House of the Hungarian Academy (the first edition was published by the latter alone). The original Hungarian version appeared in 1968 (first edition) and in 1975 (second edition).

The book is application-oriented. Numerous examples facilitate the application of the described methods.

The first two chapters give an overview on the elements of probability theory needed for mathematical statistics. The notions of sampling and order statistics are the subject of Chapter 3. This is followed by the chapters on estimation and on testing statistical hypotheses. The last two chapters of Volume 1 deal with sequential procedures and with the elements of the decision theory.

The analysis of variance, correlation and regression analysis and the statistical methods of quality control are the subject of Volume 2. Numerous tables (30 pages) facilitate the applications.

The list of references cites 115 items.

The book is clear and well presented. It can be read profitably by practitioners (engineers, chemists, physicians etc.) and mathematicians.

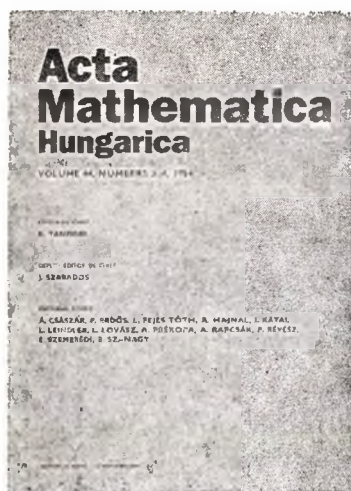
*K. Sarkadi*

# Acta Mathematica Hungarica

(Formerly: Acta Mathematica  
Academiae Scientiarum Hungaricae)

**Editor in Chief:**  
K. Tandori

**Deputy Editor in Chief:**  
J. Szabados



The journal covers a wide scope in the field of mathematics. It comprises theory of sets, mathematical logic, classical and modern analysis, algebra, number theory, geometry, topology, combinatorics, mathematical statistics, probability theory, as well as information theory.

**Founded 1950**

**Papers in English, German, French and Russian**

**Publication: two volumes annually —  
one volume contains two issues**

**Price per volume: \$ 44.00; DM 99,—**

**Size: 17 × 25 cm**

**ISSN 0236-5294**

## Order form

to be returned to

**KULTURA**

Hungarian Foreign Trading Company

P.O. Box 149, H-1389 Budapest, Hungary

- ☐ Please enter my/our subscription for  
**ACTA MATHEMATICA HUNGARICA** for one year
- ☐ Please enter my/our standing order for  
**ACTA MATHEMATICA HUNGARICA** starting with

Name: \_\_\_\_\_

Address: \_\_\_\_\_

Date and signature: \_\_\_\_\_

MACYAR  
TUDOMÁNYOS AKADÉMIA  
KÖNYVTÁRA

Contents of Volume 42. Numbers 1-2

*Ferenczi, M.*: Measures on cylindric algebras

*Petz, D.*: On spectral and central states of Banach algebras

*Györfvári, J.*: Lakunäre Interpolation mit Spline-Funktionen

*Tanović-Miller, N.*: On strong convergence of trigonometric and Fourier series

*Györy, K.*: Bounds for the solutions of norm form, discriminant form and index form equations in finitely generated integral domains

*Бейбар, К. И. и Салахова, К.*: О решетках  $N$ -радикалов, строгих слева радикалов, наследственных слева радикалов

*Grecu, E.*: Détermination des géodésiques de certains espaces riemanniens singuliers

*Parhi, N.*: On non-oscillatory solutions of second order differential inequalities

*Matolcsy, K.*: Refined extensions of syntopogenous structures and quasi-uniformities

*Alimov, Š. A. and Joó, I.*: On the eigenfunction expansions associated with the Schrödinger operator having spherically symmetrical potential

*Maknys, M.*: On the distance between consecutive prime ideal numbers in sectors

*Karamzadeh, O. A. S.*: On the Krull intersection theorem

*Kubacki, K. S. and Szynal, D.*: Weak convergence of martingales with random indices to infinitely divisible laws

*Берман, Д. Л.*: Решение одной экстремальной задачи теории операторов

*Jain, R. K.*: Semigroups with primary ideals of prime power

*Komornik, V.*: On the distribution of the eigenvalues of an orthonormal system, consisting of eigenfunctions of higher order of a linear differential operator



Akadémiai  
Kiadó

Publishing House  
of the Hungarian Academy of Sciences  
Budapest

Invitation for papers

Manuscripts should be sent to  
Acta Mathematica Hungarica  
P.O. Box 127  
H-1364 Budapest  
Hungary

PRINTED IN HUNGARY

Szegedi Nyomda, Szeged



## RECENTLY ACCEPTED PAPERS

- GUREVIČ, R., A categoricity for simplicial complices  
 ERDŐS, P., JOÓ, I. and SZÉKELY, L. A., Some remarks on infinite series  
 KELMANS, A. K., On 3-skeins in a 3-connected graph  
 JOÓ, I., On the vibration of a string  
 EWALD, G., Torische Varietäten und konvexe Polytopen  
 MOLNÁR, E., Minimal presentation of the 10 compact euclidean space forms by fundamental domains  
 BÁRÁNY, I. and KINCSES, J., Characterization of  $k$ -Helly dimensional convex bodies  
 FEJES TÓTH, G., Totally separable packing and covering with circles  
 BENT, S. W., Stable transversals and stochastic functions in polycminoes  
 BOSZNAY, A. P., A remark concerning strong uniqueness of approximations  
 HORVÁTH, M., On multidimensional universal functions  
 SARMA, M.-C. and ESCASSUT, A., Prolongement analytique à travers un T-filtre  
 PALÁSTI, I., On the seven points problem of P. Erdős  
 FEJES TÓTH, G. and HARBORTH, H., Kugelpackungen mit vorgegebenen Nachbarnzahlen  
 FENYŐ, I., On the Hankel-transformation of Schwartz distributions  
 LEFMANN, H., A note on Ramsey numbers  
 LOI, N. V. and STEINFELD, O., Gothic classes of groupoid-lattices in the theory of radicals  
 VESZTERGOMBI, K., Bounds on the number of small distances in a finite planar sets  
 DOLBILIN, N. P., Об одной характеристике решеток и неоднородной проблеме Минковского  
 RYŠKOV, S. S. and UMAROV, M. H., Строение конечных граней полиэдра  $\mu_n(m)$  при  $n \leq 4$   
 BOGNÁR, M., Walking in finite directed graphs  
 ERDŐS, P., SÁRKÖZY, A. and SÓS, V. T., Problems and results on additive properties of general sequences III  
 STROMMER, G., Über das Verhalten einer krummen Fläche in der Nähe eines parabolischen Punktes  
 FÉNYES, T., On a discrete nonlinear operational differential equation system based on the Dirichlet product  
 RACSMÁNY, A., Correction to my paper "Perfect simple Lee error-correcting codes"  
 HUSTY, M., Über eine symmetrische Schrotung mit einer Cayleyfläche als Grundfläche  
 LOI, N. V., A note on the radical theory of involution algebras  
 DORNINGER, D. and LÄNGER, H., On a set of relations arising from the triangulation problem  
 ABRAMS, G. D., The recovery question for local incidence rings  
 SZABADOS, J., VARMA, A. K. and SELVARAJ, C. R., Error estimates of a general lacunary trigonometric interpolation of equidistant nodes  
 RÖSCHEL, O., Torusflächen des Galileischen Raumes  $G_3$   
 PIOCHI, B., Congruences on inversive hemirings  
 VELDSMAN, S., Hereditary conditions on classes of near-rings  
 HORVÁTH, M., Answer to a problem of I. Joó  
 BOOTH, G. L., A note on  $\Gamma$ -near-rings  
 YUSUF, S. M. and SHABIR, M., Radical classes and semisimple classes for hemirings  
 TOSIĆ, R., On cops and robber game  
 CASTRO, S., Miquelsche Minkowski-Ebenen in spiegellungsgeometrischer Darstellung

## CONTENTS (continued)

KHAN, L. A., Separability in the uniform topology .....	407
EIGEN, S. J., Putting convergent sequences into measurable sets .....	411
SIMON, L., On approximation of solutions of exterior boundary value problems .....	413
DAMASCHKE, P. and STERN, M., A characterization of generalized matroid lattices .....	425
FIALOWSKI, A., On the deformations of $L_1$ .....	433
REIMNITZ, P., An arcsine-law for the oscillating random walk .....	439
DUGGAL, B. P., On the spectrum of a class of integral transforms II .....	451
DUGGAL, B. P., On the spectrum of a class of integral transforms III. An application .....	461
FLORIAN, A., On compact packing of circles .....	473
BOOK REVIEWS .....	481

# CONTENTS

IMHOF, J. P., On Brownian bridge and excursion .....	1
BIHARI, I., An asymptotic statement concerning the solutions of the differential equation $x'' + a(t)x = 0$ .....	11
BIHARI, I., Note to an extension of a Sturmian comparison theorem .....	15
BELL, H. E., On commutativity of quasi-commutative rings .....	21
HUISMANS, C. B., An inequality in complex Riesz algebras .....	29
FRANKL, P., Bounding the size of a family knowing the cardinality of differences .....	33
KHOI, TRINH DANG, Строго наследственные радикалы в классе всех топологических колец .....	37
BEAZER, R., Congruence uniform algebras with pseudocomplementation .....	43
BLASCO, J. L., Complete bases in topological spaces .....	49
GROSSMAN, E. H., Number bases in quadratic fields .....	55
HUYNH, DINH VAN, On rings with modified chain conditions .....	59
GUT, A., On the law of the iterated logarithm for randomly indexed partial sums with two applications .....	63
STEIN, S., Lattice-tiling by certain star bodies .....	71
SUBBARAO, M. V. and SITARAMACHANDRARAO, R., The distribution of values of a class of arithmetic functions .....	77
POPENDA, J., On the boundedness of the solutions of second order differential equations ....	89
Дарбинян, С. X., Панцикличность орграфов при условии Мейнила .....	95
PANNY, W. and PRODINGER, H., The expected height of paths for several notions of height .....	119
AHMAD, M., Estimation of the parameters of Burr distribution based on order statistics .....	133
COLBOURN, C. J. and ROSA, A., Indecomposable triple systems with $\lambda=4$ .....	139
LÉNÁRD, M., Spline interpolation in two variables .....	145
ISAC, G., Branches continues de vecteurs propres généralisés. Applications aux équations de coïncidences .....	155
XEKALAKI, E., Some bivariate extensions of the generalized Waring distribution .....	173
IVIĆ, A., The distribution of primitive abundant numbers .....	183
IVIĆ, A., On squarefree numbers with restricted prime factors .....	189
BERG, G., Steinness and the vanishing of cohomology .....	193
SEOH, M. and PURI, M. L., Berry—Esséen theorems for signed linear rank statistics under near location alternatives .....	197
KHARE, S. S., Reduction of equivariant bordism groups .....	213
FEJES TÓTH, L., Packing of homothetic discs of $n$ different sizes .....	217
SRIVASTAVA, K. B., A remark on Mathur's paper. Simple proof of Telyakovskii—Gopengauz's theorem .....	223
CHIANG, C.-Y. and PURI, M. L., Tests of subhypotheses in linear regression based on rank-order estimates .....	237
BECK, J., Remarks on combinatorial geometry I .....	249
VEIDINGER, L., On the order of convergence of a finite element method for the biharmonic equation .....	255
FLORIAN, A. and GROEMER, H., Two remarks on the permeability of layers of convex bodies .....	259
MARCUS, F., Sur les surfaces à groupes continus $G_2$ de similitude projectives en elles-mêmes et sur les surfaces complexes .....	267
KOMORNIK, V., On the eigenfunctions of first- and second-order differential operators .....	275
KISS, P., Differences of the terms of linear recurrences .....	285
HARMAN, G., PINTZ, J. and WOLKE, D., A note on the Möbius and Liouville functions .....	295
PLONKA, J., On the sum of a $\iota$ -semilattice ordered system of algebras .....	301
MILLER, H. I. and XENIKAKIS, P. J., Some results on the Cantor set .....	309
Баясгалан, Ц., О фундаментальной приводимости самосопряженных и унитарных операторов в пространствах с индефинитной метрикой .....	313
POYATOS, F., Archimedean decompositions of left $S$ -semimodules and semirings .....	323
GRIMMETT, G. R., The largest components in a random lattice .....	325
GAÁL, I., Norm form equations with several dominating variables and explicit lower bounds for inhomogeneous linear forms with algebraic coefficients II .....	333
PERELLI, A. and SALERNO, S., On $2k$ -dimensional density estimates .....	345
SALERNO, S., On the distribution of $x_1^2 + \dots + x_n^2$ in the arithmetic progressions .....	357
LÖFSTRÖM, J., Best approximation in $L_p(w)$ by algebraic polynomials .....	375
MOÖR, A. and NADJ, D. F., Über die autoparallele Abweichung von Finsler—Otsukischen Räumen und Anwendungen in Räumen mit speziellen $P$ -Tensoren .....	395

(continued inside)